

The Ocean Colour Climate Change Initiative: III. A round-robin comparison on in-water bio-optical algorithms

Robert J.W. Brewin^{*,a,b}, Shubha Sathyendranath^{a,b}, Dagmar Müller^c, Carsten Brockmann^d, Pierre-Yves Deschamps^e, Emmanuel Devred^f, Roland Doerffer^c, Norman Fomferra^d, Bryan Franz^g, Mike Grant^a, Steve Groom^a, Andrew Horseman^a, Chuanmin Hu^h, Hajo Krasemann^c, ZhongPing Leeⁱ, Stéphane Maritorena^j, Frédéric Mélin^k, Marco Peters^d, Trevor Platt^a, Peter Regner^l, Tim Smyth^a, Francois Steinmetz^e, John Swinton^m, Jeremy Werdell^g, George N. White IIIⁿ

^aPlymouth Marine Laboratory (PML), Prospect Place, The Hoe, Plymouth PL1 3DH, UK

^bNational Centre for Earth Observation, PML, Plymouth PL1 3DH, UK

^cHelmholtz-Zentrum Geesthacht, Max-Planck-Straße 1, 21502 Geesthacht, Germany

^dBrockmann Consult, Max-Planck-Straße 2, D-21502 Geesthacht, Germany

^eHYGEOS, 165 Avenue de Bretagne, 59000 Lille, France

^fUniversité Laval, 2325, rue de l'Université, Québec G1V 0A6, Canada

^gNASA Goddard Space Flight Center, Greenbelt, Maryland, USA

^hCollege of Marine Science, University of South Florida, St. Petersburg, FL 33701, USA

ⁱCollege of Science and Mathematics, University of Massachusetts Boston, Boston, MA 02125-3393, USA

^jUniversity of California Santa Barbara, Santa Barbara, CA 93106, USA

^kEuropean Commission - Joint Research Centre, Institute for Environment and Sustainability, Ispra, 21027, Italy

^lEuropean Space Agency, ESRIN, Via Galileo Galilei, Casella Postale 64 00044 Frascati, Italy

^mTelespazio VEGA UK Ltd, 350 Capability Green, Luton, Bedfordshire, LU1 3LU, UK

ⁿOcean Science Division, Bedford Institute of Oceanography, Box 1006, Dartmouth, Nova Scotia, B2Y 4A2, Canada

Abstract

Satellite-derived remote-sensing reflectance (R_{rs}) just above the sea surface can be used for mapping biogeochemically relevant variables, such as the chlorophyll concentration and the Inherent Optical Properties (IOPs) of the water, at global scales for use in climate-change studies. Prior to generating such prod-

ucts, suitable algorithms have to be selected that are appropriate for the purpose. Algorithm selection needs to account for both qualitative and quantitative requirements. In this paper, we develop an objective methodology designed to rank the quantitative performance of a suite of bio-optical models. The objective classification is applied using the NASA bio-Optical Marine Algorithm Data set (NOMAD). Using *in situ* R_{rs} as input to the models, the performance of eleven semi-analytical models, as well as five empirical chlorophyll algorithms and an empirical diffuse attenuation coefficient algorithm, are ranked for spectrally-resolved IOPs, chlorophyll concentration and the diffuse attenuation coefficient at 489 nm. The sensitivity of the objective classification and the uncertainty in the ranking is tested using a Monte-Carlo approach (bootstrapping). Results indicate that the performance of the semi-analytical models varies depending on the product and wavelength of interest. For chlorophyll retrieval, empirical algorithms perform better than semi-analytical models, in general. The performance of these empirical models reflect either their immunity to scale errors or instrument noise in R_{rs} data, or simply that data used for model parameterisation were not independent of NOMAD. Nonetheless, uncertainty in the classification suggest the performance of some semi-analytical algorithms at retrieving chlorophyll were comparable with the empirical algorithms. For phytoplankton absorption at 443 nm, some semi-analytical models also performed with similar accuracy to an empirical model. We discuss the potential biases, limitations and uncertainty in the approach, as well as additional qualitative considerations for algorithm selection for climate change studies. Our classification has the

potential to be routinely implemented, such that the performance of emerging algorithms can be compared with existing algorithms as they become available. In the long-term, such an approach will further aid algorithm development for ocean-colour studies.

29 *Key words:* Phytoplankton, Ocean colour, Inherent Optical Properties, Remote
30 sensing, chlorophyll-a

31 **1. Introduction**

32 Visible radiance received by satellite ocean-colour sensors over oceanic re-
33 gions is essentially influenced by two components: the atmosphere and the
34 ocean. Typically, the atmospheric component constitutes more than 80% of
35 the signal received by the sensor, and it needs to be removed to isolate the
36 signal from the ocean. The ocean-colour signal may then be used to quan-
37 tify optically-significant water-constituents such as Coloured Dissolved Organic
38 Matter (CDOM) and the abundance of particulate matter, inclusive of phyto-
39 plankton, indexed through their chlorophyll pigment concentration, and non-
40 phytoplanktonic material (e.g. detrital and inorganic matter).

41 Phytoplankton are a key component of the Earth System and are recognised
42 as an Essential Climate Variable in the Implementation Plan of the Global Cli-
43 mate Observing System (GCOS, 2011). Phytoplankton absorb light energy that
44 is either dissipated as heat, directly influencing the physical properties of the

*Corresponding author. Plymouth Marine Laboratory (PML), Prospect Place, The Hoe, Ply-
mouth PL1 3DH, UK

Email address: `robr@pml.ac.uk` (Robert J.W. Brewin)

45 oceans, or used for photosynthesis (primary production), by which light is con-
46 verted into chemical energy and carbon converted from inorganic to organic
47 form. It is estimated that phytoplankton fix approximately 50 gigatons of car-
48 bon per year, equivalent to net terrestrial primary production. Phytoplankton,
49 together with physical processes, regulate the CO₂ concentration of the surface
50 ocean and the rate of CO₂ exchanges between the atmosphere and ocean. They
51 are at the base of the food web, providing sustenance for all pelagic marine life,
52 and contribute to the biogeochemical cycling of a variety of climatically-relevant
53 elements, such as silica, nitrate and phosphate, in addition to carbon. Monitor-
54 ing the variability in phytoplankton distribution is vital to understanding how the
55 ocean ecosystem is likely to respond to future changes in climate.

56 The concentration of CDOM, its photodegradation status and the concentra-
57 tion of detrital matter present in the water have a significant effect on phyto-
58 plankton photosynthesis, through their absorption of light at blue wavelengths
59 of the visible spectrum, which corresponds to the main phytoplankton absorp-
60 tion peak. CDOM can also affect the transport and bioavailability of trace metals
61 (Santschi et al., 1997; Guo et al., 2001), with possible implications for biological
62 activity, and plays an important role in photochemistry and photobiology, with
63 implications for ocean-climate connections (Nelson and Siegel, 2013). The pres-
64 ence of highly-scattering non-phytoplanktonic particulate material (e.g. detrital
65 and inorganic matter) alters the spectral quality of the underwater light field and
66 thus influences phytoplankton photosynthesis. The concentration of particulate
67 material in the water is also important in coastal regions and has implications for

68 coastal protection, shipping and recreational activities. These are some of the
69 reasons why the systematic monitoring of ocean colour is considered a require-
70 ment for climate research by GCOS (GCOS, 2011) and why it is a component of
71 the Climate Change Initiative (CCI) of the European Space Agency (ESA).

72 The CCI programme was launched to realise the full potential of long-term,
73 global, Earth Observation archives that ESA as well as its member states have
74 established over the past 30-years, and to contribute to the Essential Climate
75 Variable databases required by United Nations Framework Convention on Cli-
76 mate Change (UNFCCC). The Ocean Colour CCI (OC-CCI) project is one of 14
77 ESA funded CCI projects. The aims of OC-CCI are to create a long-term, consis-
78 tent, error-characterised time series of ocean-colour products, for use in climate
79 change studies. A key component of the programme is the selection of suitable
80 algorithms that meet user requirements and project aims. The selection of algo-
81 rithms for the OC-CCI project can be partitioned into two parts: (i) selection of
82 algorithms that correct for atmospheric affects; and (ii) algorithms that convert
83 the retrieved ocean-colour signal into biogeochemically relevant variables, here-
84 after referred to as atmospheric-correction and in-water algorithms respectively.
85 This paper focuses on the development of an objective methodology designed to
86 aid the selection of appropriate in-water algorithms for climate studies. For infor-
87 mation regarding the selection of atmospheric-correction algorithms the reader
88 is referred to Müller et al. (Submitted) in this issue.

89 Since the establishment of ocean-colour remote sensing from space, with
90 the launch of the Coastal Zone Color Scanner (CZCS) of NASA on board the

91 Nimbus-7 satellite in 1978, blue-to-green ratios of water-reflectance have been
92 used in empirical relationships to derive the total concentration of chlorophyll-
93 a (C), an ubiquitous pigment present in phytoplankton. With the launch of
94 the Sea-viewing Wide Field-of-view Sensor (SeaWiFS), the NASA successor
95 to CZCS, NASA organised the SeaWiFS Bio-optical Algorithm Mini-workshop
96 (SeaBAM; O'Reilly et al., 1998), designed to identify chlorophyll algorithms
97 suitable for operational use for processing SeaWiFS data. A database was devel-
98 oped with simultaneous measurements of *in situ* chlorophyll and *in situ* measure-
99 ments of remote-sensing reflectance just above the surface ($R_{rs}(\lambda)$). Based on the
100 results from the workshop, an empirical blue-green band ratio algorithm, labelled
101 the Ocean-Chlorophyll-2 (OC2) algorithm, was chosen as the operational algo-
102 rithm for SeaWiFS. This was later updated to the Ocean-Chlorophyll-4 (OC4)
103 algorithm (O'Reilly et al., 2000).

104 In Case-1 waters (Morel and Prieur, 1977) typically encountered in the open
105 ocean, where variations in ocean-colour are driven primarily by the abundance of
106 phytoplankton, with a co-varying influence from particulate matter and CDOM,
107 empirical blue-green band-ratio algorithms were generally found to perform with
108 reasonable accuracy. However, in more optically-complex waters (Case-2 wa-
109 ters according to Morel and Prieur, 1977), often encountered in coastal regions,
110 where the concentrations of particulate matter and CDOM do not covary in a
111 predictable manner with the abundance of phytoplankton, empirical blue-green
112 band-ratio algorithms can give spurious results (e.g. Lavender et al., 2004).

113 Theoretical approaches have demonstrated that $R_{rs}(\lambda)$ is related to the In-

114 herent Optical Properties (IOPs) of seawater, the absorption and backscattering
115 coefficients. The absorption coefficient can in turn be partitioned into the contri-
116 butions from water itself, and the type and abundance of material present in the
117 water, including phytoplankton, detrital matter and CDOM. The backscattering
118 coefficient can be partitioned into contributions from pure seawater and partic-
119 ulate matter suspended in the water (which includes phytoplankton). IOPs can
120 be used to infer biogeochemical processes and to estimate the concentrations of
121 various optically-significant water constituents, such as chlorophyll. Theoreti-
122 cal approaches that derive IOPs from $R_{rs}(\lambda)$ may improve performance of algo-
123 rithms in more optically-complex waters (see IOCCG, 2000), and a variety of
124 semi-analytical approaches have been developed in this direction (see IOCCG,
125 2006).

126 Recently, NASA organised an international IOP algorithm workshop (Werdell,
127 2009), designed to provide data sets (Werdell and Bailey, 2005) and processing
128 framework in an international forum within which a new generation of global
129 IOP products can be developed and evaluated. The workshop aimed to: define
130 the state of the art with regard to the application of semi-analytical models to
131 satellite radiometry; identify similarities and differences between approaches;
132 identify strategies to provide uncertainties in IOPs; and achieve community con-
133 sensus toward the generation of global IOP products (Werdell, 2009). An output
134 of the workshop was the development of a Generalised Inherent Optical Prop-
135 erty model (GIOP), a test platform for algorithm development that offers free-
136 dom to specify various optimisation approaches and parameterisations (Franz

137 and Werdell, 2010; Werdell et al., 2013).

138 In contrast to the aims of the NASA GIOP workshop, but making use of
139 progress made as a result of the workshop, and building on the report of the
140 IOCCG working group on the topic (IOCCG, 2006), this paper aims to establish
141 an objective methodology for algorithm selection for climate-change studies, and
142 then to use the method to compare and rank a variety of algorithms. Both qual-
143 itative and quantitative considerations are examined. Qualitative considerations
144 relate to the suitability of the algorithms for use in climate change studies and
145 the quantitative considerations relate to algorithm performance. Qualitative al-
146 gorithm considerations include the ability of the algorithm to:

- 147 • Create a long-term, consistent, error-characterised time series of ocean-
148 colour products for use in climate-change studies;
- 149 • Generate products that best suit the requirements of the user community;
- 150 • Facilitate seamless merging of Case-1 (open-ocean) and Case-2 (coastal
151 optically-complex) waters;
- 152 • Quantify a variety of properties of the marine ecosystem that are relevant
153 to climate studies and accessible from satellite ocean-colour data and;
- 154 • Be robust against potential modifications in the marine ecosystem in a
155 changing climate.

156 Ideally, the most suitable algorithm would meet all these requirements and com-
157 pare well in statistical tests of performance. Using a suite of statistical tests,

158 and an *in situ* database of chlorophyll (C), the diffuse attenuation coefficient at
159 489 nm ($K_d(489)$), IOPs and $R_{rs}(\lambda)$, we evaluate the quantitative performance
160 of a number of empirical and semi-analytical in-water bio-optical models. The
161 limitations of the approach are discussed and additional challenges regarding the
162 selection of in-water algorithms for climate studies are highlighted.

163 2. Data

164 To test in-water bio-opticals models, we made use of the publicly-available
165 NASA bio-Optical Marine Algorithm Data set (NOMAD, Werdell and Bailey,
166 2005). NOMAD Version 2.0 ALPHA was compiled on 18 July 2008 by the
167 NASA Ocean Biology Processing Group and source data is available online
168 (<http://seabass.gsfc.nasa.gov/seabasscgi/nomad.cgi>), as is documentation related
169 to IOPs (Werdell, 2005). The NOMAD database provides global *in situ* measure-
170 ments of above-water spectral water-leaving radiance ($L_w(\lambda)$) and spectral sur-
171 face irradiance ($E_s(\lambda)$), from which remote-sensing reflectance can be computed
172 ($R_{rs}(\lambda) = L_w(\lambda)/E_s(\lambda)$), and coincident measurements of water constituents such
173 as the chlorophyll-a concentration, IOPs and $K_d(489)$ (diffuse attenuation coef-
174 ficient at 489 nm). The solar sun-zenith angle (θ) was computed for each data
175 point using information on time and location. Table 1 denotes the variables used
176 in the comparison.

177 The OC-CCI project currently focuses on the use of three ocean-colour satel-
178 lite platforms: the Medium Resolution Imaging Spectrometer (MERIS) of ESA;
179 the Moderate Resolution Imaging Spectro-radiometer (MODIS) of NASA; and
180 the Sea-viewing Wide Field-of-view Sensor (SeaWiFS) of NASA, to create a

time-series of satellite data. Therefore, to be representative of the majority of wavelengths in all three satellite sensors, a common band set of 411, 443, 489, 510, 555, and 665 nm was chosen to maximise the amount of validation data points in NOMAD. Though there are some mis-matches (MERIS native 560 > 555 nm; MODIS native 547 < 555 nm and 531 nm excluded; and SeaWiFS 670 > 665 nm), this compromise was adopted to maximise the number of samples. The common band set used included six bands compatible with MERIS and SeaWiFS and five bands compatible with MODIS. Co-located *in situ* measurements of $R_{rs}(\lambda)$ were used as input to the models, as opposed to satellite-derived $R_{rs}(\lambda)$, to minimise mis-matches in spatial scales between input and output variables.

To maximise the number of $b_b(\lambda)$ samples, 670 nm was used where reflectance data at 665 nm were unavailable. Note that $b_b(\lambda)$, and the slope of $b_b(\lambda)$, denoted as γ (Table 1), were used in this comparison as opposed to partitioning $b_b(\lambda)$ into the contribution from pure water (b_{bw}) and particles (b_{bp}), to avoid issues caused by different b_{bw} spectra in different semi-analytical models. Remote sensing reflectance data, at various wavelengths, and solar-zenith angles were used as input to in-water algorithms to estimate IOPs, C and $K_d(489)$ (Table 1, 2 and 3). Estimated variables using the models were then compared with *in situ* values in NOMAD, to determine the performance of the algorithms. Figure 1 shows the spatial coverage and number of samples for each variable used in the *in situ* database and the NOMAD record identifier for each measurement used in the comparison is provided as Supplementary Data.

203 3. Models

204 The following sections describe the semi-analytical models, designed to re-
205 trieve IOPs, and the chlorophyll models and the diffuse attenuation coefficient
206 (K_d) models incorporated into the comparison. Tables 2 and 3 also provide a
207 description of the output variables of each model and a summary listing key at-
208 tributes of the various algorithms.

209 3.1. Semi-analytical models

210 Semi-analytical models used in the comparison are described in this section.
211 The term ‘semi-analytical models’ will be conventionally employed hereafter to
212 describe Models A-K for the sake of brevity. However, we acknowledge that
213 some of the models vary in their use of analytical and empirical solutions to
214 solve for the IOPs. These semi-analytical models (A-K) are used to compute the
215 total absorption coefficient (a), combined absorption by detritus and coloured
216 dissolved organic matter or gelbsfoff (a_{dg}), absorption by phytoplankton (a_{ph}),
217 total back-scattering coefficient (b_b), the spectral slope of the total backscatter-
218 ing coefficient (γ), the spectral slope of a_{dg} , denoted S_{dg} , and the ratio of phyto-
219 plankton absorption at 555 nm to that at 443 nm ($a_{ph}(555)/a_{ph}(443)$) (see Table 1
220 for all notations used). The ratio $a_{ph}(555)/a_{ph}(443)$ was used in this comparison
221 as an index of the spectral shape of the phytoplankton absorption coefficient, an
222 index of the community structure of the phytoplankton (Sathyendranath et al.,
223 2001, 2004; Ciotti et al., 2002). The ratio of 555 nm to 443 nm was chosen as
224 these wavelengths typically represent the minimum and maximum of the phy-
225 toplankton absorption spectra. However, we acknowledge that ratios of other

226 wavelengths could have also been used.

227 3.1.1. Model A

228 Model A refers to the model of Smyth et al. (2006). It uses an algebraic
229 approach for determining IOPs. The model uses spectral slopes for $a - a_w$ (where
230 subscript w stands for water) and b_{bp} (total particulate backscattering) derived
231 from field measurements, at the central wavelengths of 490 and 510 nm (or 531
232 for MODIS). Once the absorption and backscattering coefficients are known at
233 these wavelengths, based on Morel (1980), and assuming a fixed spectral slope
234 for b_{bp} , the absorption and backscattering coefficients across the spectrum can be
235 determined. Once absorption and backscattering are determined spectrally, a_{dg}
236 and a_{ph} can be determined using standard relationships and slopes between the
237 wavelengths of 412 and 443 nm.

238 3.1.2. Model B

239 Model B refers to the model of Smyth et al. (2006), as in Model A, but apply-
240 ing a new optical water classification, whereby the model parameters (spectral
241 slopes in $a - a_w$ and b_{bp}) were computed for eight optical classes (see Moore
242 et al., 2009). Based on the fuzzy-class-membership for each sample, determined
243 from R_{rs} , the spectral slopes in $a - a_w$ and b_{bp} are re-computed and implemented
244 in the model of Smyth et al. (2006): a_{dg} and a_{ph} are then determined as in Model
245 A.

246 3.1.3. Model C

247 Model C refers to the ocean-colour model of Devred et al. (2011) with some
248 simplifications. This model is designed to derive in-water optical properties and
249 water constituents from spectral water-leaving radiances, using non-linear op-
250 timisation procedures. The method makes use of a three-component model of
251 phytoplankton absorption coupled to the reflectance model of Sathyendranath
252 and Platt (1997). The model retrieves $b_{bp}(555)$ (assuming the slope of $b_{bp} = 1.03$
253 following Maritorena et al. (2002)), $a_{dg}(443)$ and S_{dg} from R_{rs} , initially assum-
254 ing that a_{ph} can be expressed as the sum of the absorption coefficient of three
255 phytoplankton size classes (pico-, nano- and micro-phytoplankton), each with
256 its particular specific absorption spectrum (a_{ph}^* , phytoplankton absorption nor-
257 malised by chlorophyll concentration) derived from the NOMAD dataset. Wave-
258 lengths from 443 to 555 nm were used in the inversion of Model C. Output vari-
259 ables were constrained to lie within the following range: $0.0 < a_{ph} < 100 \text{ m}^{-1}$;
260 $0.0 < a_{dg} < 100 \text{ m}^{-1}$; and $0.0 < b_{bp} < 5.0 \text{ m}^{-1}$.

261 3.1.4. Model D

262 Model D refers to the algebraic Quasi-Analytical Algorithm (QAA) of Lee
263 et al. (2002). The model was designed to retrieve IOPs in optically-deep waters.
264 The model inversion is based on two steps: the first involves partitioning water
265 reflectance into b_b and a and the second decomposing a into a_{dg} and a_{ph} . The
266 model is referred to as “Quasi-Analytical” as parts of the inversion are based
267 on analytical, semi-analytical and empirical approximations. Model D uses the
268 original parameterisation as described in Lee et al. (2002).

269 *3.1.5. Model E*

270 Model E refers to the Quasi-Analytical Algorithm (QAA) of Lee et al.
271 (2002), as in Model D, but following an updated parameterisation (see Lee et al.,
272 2009). This includes the use of measured $R_{rs}(670)$ in the calculation of $a(555)$,
273 in contrast to Model D which instead uses $R_{rs}(640)$ in the calculation of $a(555)$,
274 estimated empirically from other wavelengths when using data from SeaWiFS,
275 MODIS, or MERIS.

276 *3.1.6. Model F*

277 Model F refers to the physics-based Hyperspectral Optimization Process
278 Exemplar (HOPE) model of Lee et al. (1998, 1999). In this model, R_{rs} is
279 modelled as a function of IOPs, and when influencing the R_{rs} signal, bottom
280 depth and bottom albedo. Unknowns are derived from non-linear optimisa-
281 tion. The spectral shape of bottom albedo is pre-determined before the opti-
282 misation starts, with the choice of two shapes (one for sand, another for grass)
283 automatically selected using the R_{rs} spectrum. The phytoplankton absorption
284 coefficients were constrained to lie within an upper and lower boundary (e.g.
285 $0.002 < a_{ph}(443) < 1.0 \text{ m}^{-1}$).

286 *3.1.7. Model G*

287 Model G refers to the semi-analytical Garver-Siegel-Maritorena (GSM)
288 model, that was initially developed by Garver and Siegel (1997) and later up-
289 dated by Maritorena et al. (2002). The GSM model retrieves simultaneous esti-
290 mates of chlorophyll (C), $a_{dg}(443)$ and $b_{bp}(443)$ from $R_{rs}(\lambda)$, assuming an under-
291 lying bio-optical model and using non-linear optimisation. Global parameters

292 of the bio-optical model were initially assigned based on simulated annealing
 293 on a global quasi-real dataset, which are then used in the non-linear optimisa-
 294 tion routine. These include a fixed chlorophyll-specific phytoplankton absorp-
 295 tion coefficient (a_{ph}^*), S_{dg} and the slope of b_{bp} . The chlorophyll (C), $a_{dg}(443)$
 296 and $b_{bp}(443)$ are first retrieved by fitting the bio-optical model to the observed
 297 $R_{rs}(\lambda)$. IOPs at any wavelengths are then obtained using C , $a_{dg}(443)$ and $b_{bp}(443)$
 298 and their specific shape function from the bio-optical model. For Model G, the
 299 output variables are constrained to lie within the range that was used to param-
 300 eterise the model ($0.01 < C < 64 \text{ mg m}^{-3}$; $0.0001 < a_{dg}(443) < 2.0 \text{ m}^{-1}$; and
 301 $0.0001 < b_{bp}(443) < 0.1 \text{ m}^{-1}$).

302 3.1.8. Model H

303 Model H refers to the semi-analytical Garver-Siegel-Maritorena (GSM)
 304 model (Maritorena et al., 2002), as in Model G, but allowing the retrievals to
 305 have any value, thus removing the constraint imposed on Model G.

306 3.1.9. Model I

307 Model I refers to a preliminary configuration of the Generalized Inherent Op-
 308 tical Property algorithm (GIOP; Franz and Werdell, 2010; Werdell et al., 2013).
 309 The GIOP model is designed as a test platform for algorithm development and
 310 was the result of a NASA IOP Algorithm Workshop (see Werdell, 2009; Werdell
 311 et al., 2013). Whereas the GIOP model offers the user freedom to specify differ-
 312 ent parameterisations and optimisation approaches, a preliminary configuration
 313 for GIOP is available which includes: an assigned a_{ph}^* following Bricaud et al.
 314 (1995) but normalised by $0.055 \text{ m}^2 (\text{mgC})^{-1}$; a spectral backscattering depen-

315 dency following the QAA; a fixed spectral slope for $a_{dg}(\lambda)$ of 0.018 nm^{-1} ; Morel
 316 et al. (2002) f/Q ratio for zero Sun angle and zero view angle, where $Q(\lambda)$ is the
 317 ratio of upwelling irradiance to upwelling radiance and $f(\lambda)$ captures the net ef-
 318 fects of variation in sea state, illumination conditions, and water column content;
 319 and Levenberg-Marquardt optimisation. It is designed to retrieve spectral IOPs
 320 and chlorophyll, and it is worth noting that this preliminary configuration could
 321 be changed with time. All IOPs (a_{dg} , a_{ph} , b_{bp} , and $a_{dg} + a_{ph}$) were constrained
 322 to lie within -0.005 and 5 m^{-1} . Retrievals were excluded if the reconstructed
 323 R_{rs} spectrum, between 411-555 nm, differed from the observed R_{rs} spectrum by
 324 more than 33%.

325 3.1.10. Model J

326 Model J refers to a Case-1 model, in which all IOPs are modelled as a func-
 327 tion of the chlorophyll concentration (C) derived using the NASA OC4v6 em-
 328 pirical model (Model L). Once C is estimated from R_{rs} , C is used as input to
 329 estimate: $a_{ph}(\lambda)$ using a three-component model of phytoplankton absorption
 330 (Brewin et al., 2011); $a_g(\lambda)$ using a power-function of C (Morel, 2009) with an
 331 exponential spectral slope (S_g) of 0.018 nm^{-1} ; $a_d(\lambda)$ using a power-function of
 332 C (Bricaud et al., 2010) with an exponential spectral slope (S_d) of 0.0094 nm^{-1} ;
 333 $b_{bp}(\lambda)$ as a function of C using the model of Huot et al. (2008); pure water
 334 absorption (a_w) according to Pope and Fry (1997); and pure-water backscatter-
 335 ing (b_{bw}) according to Buiteveld et al. (1994). Components of absorption and
 336 backscattering are added to obtain the totals a and b_b respectively, from which
 337 R_{rs} is computed using the model of Gordon et al. (1988).

338 3.1.11. Model K

339 Model K refers to a preliminary configuration of an in-water artificial
340 Neural-Network (NN) (e.g. Doerffer and Schiller, 2000, 2006; Doerffer et al.,
341 2002) which is used as the forward model within an optimisation procedure
342 (Levenberg-Marquardt). The model computes IOPs from water-leaving radiance
343 for all available multi-spectral ocean colour sensors as well as *in situ* measure-
344 ments. The method was optimised to invert water-leaving radiance directly into
345 spectral IOPs, with chlorophyll (*C*) parameterised as a function of phytoplankton
346 absorption and $K_d(489)$ as a function of scattering and total absorption.

347 3.2. Chlorophyll (*C*) models

348 Chlorophyll (*C*) algorithms incorporated into the comparison are described
349 in the following section. For semi-analytical Models C, G, H, I, and K, chloro-
350 phyll is an output from the models. For semi-analytical Models A, B, D, E, and
351 F, chlorophyll is not an output. For the purposes of the comparison, we esti-
352 mated chlorophyll as a function of $a_{ph}(443)$ using a power-law relationship (e.g.
353 Bricaud et al., 1995), such that

$$C = \left[\frac{a_{ph}(443)}{A} \right]^{\frac{1}{B}}, \quad (1)$$

354 where, A and B are positive empirical parameters. The empirical parameters A
355 and B were computed using the *in situ* NOMAD database (1042 samples), and
356 set to $A = 0.0497$ and $B = 0.7575$. For semi-analytical Models A, B, D, E, and
357 F, $a_{ph}(443)$ was first computed, then chlorophyll was computed using Eq. (1).

358 It is worth noting that the empirical conversion from $a_{ph}(443)$ to chlorophyll is
 359 merely introduced to facilitate the comparison, it is not a feature of the original
 360 algorithms. Note that Model J is not incorporated in the chlorophyll comparison
 361 as this model uses chlorophyll estimated from an empirical model (Model L)
 362 as input to compute IOPs. In addition to the semi-analytical models (A-I and
 363 K), a variety of empirical chlorophyll algorithms were also incorporated into the
 364 comparison and are described below.

365 3.2.1. Model L

366 Model L refers to the NASA OC4 chlorophyll algorithm (O'Reilly et al.,
 367 2000). This is a polynomial algorithm that relates the log-transformed ratio
 368 of remote-sensing reflectances (X) to the chlorophyll concentration (C). The
 369 OC4v6 uses a four-band blue-green reflectance ratio such that:

$$X = \log_{10}\{[R_{rs}(443) > R_{rs}(489) > R_{rs}(510)]/R_{rs}(555)\}. \quad (2)$$

370 Chlorophyll (C) is estimated according to:

$$C = 10^{(a_0 + a_1 X + a_2 X^2 + a_3 X^3 + a_4 X^4)}, \quad (3)$$

371 where, $a_0 = 0.3272$, $a_1 = -2.9940$, $a_2 = 2.7218$, $a_3 = -1.2259$ and $a_4 =$
 372 -0.5683 (NASA, 2010).

373 3.2.2. Model M

374 Model M refers to the NASA OC3S chlorophyll algorithm (O'Reilly et al.,
375 2000). Like the OC4, this is a polynomial algorithm that relates the log-
376 transformed ratio of remote-sensing reflectances (X) to the chlorophyll concen-
377 tration (C). The OC3S uses a three-band blue-green reflectance ratio where

$$X = \log_{10}\{[R_{rs}(443) > R_{rs}(489)]/R_{rs}(555)\}, \quad (4)$$

378 and chlorophyll (C) is estimated according to Eq. (3) where, $a_0 = 0.2515$, $a_1 =$
379 -2.3798 , $a_2 = 1.5823$, $a_3 = -0.6372$ and $a_4 = -0.5692$ (NASA, 2010).

380 3.2.3. Model N

381 Model N refers to the NASA OC2S chlorophyll algorithm (O'Reilly et al.,
382 2000). Like the OC4 and OC3S, this is a polynomial algorithm that relates the
383 log-transformed ratio of remote-sensing reflectances (X) to the chlorophyll con-
384 centration (C). The OC2S uses a two-band blue-green reflectance ratio where

$$X = \log_{10}[R_{rs}(489)/R_{rs}(555)], \quad (5)$$

385 and chlorophyll (C) is estimated according to Eq. (3) where, $a_0 = 0.2511$, $a_1 =$
386 -2.0853 , $a_2 = 1.5035$, $a_3 = -3.1747$ and $a_4 = 0.3383$ (NASA, 2010).

387 3.2.4. Model O

388 Model O refers to the MERIS chlorophyll band-ratio algorithm (Morel and
389 Antoine, 2011). Like the OC4, it is a four-band polynomial algorithm that relates

the log-transformed ratio of remote-sensing reflectance (X) to the chlorophyll concentration (C). Considering that a common-band set was chosen, not inclusive of 560 nm, the algorithm was implemented following Morel et al. (2007), such that the wavelength of 560 nm was replaced by 555 nm, and X can be estimated following Eq. (2) and (3), where $a_0 = 0.4461529$, $a_1 = -3.291807$, $a_2 = 3.777216$, $a_3 = -4.172339$ and $a_4 = 1.415588$ (see Table 2 OC4Me555 of Morel et al., 2007).

3.2.5. Model P

Model P refers to the chlorophyll algorithm of Hu et al. (2012). This empirical algorithm was designed to improve the estimate of chlorophyll (C) in the global ocean at concentrations $\leq 0.25 \text{ mg m}^{-3}$. For low chlorophyll concentrations ($\leq 0.25 \text{ mg m}^{-3}$), the algorithm uses a colour index (CI), which is defined as the difference between R_{rs} in the green region of the visible spectrum and a reference formed linearly between R_{rs} in the blue and red region of the visible spectrum. For high chlorophyll concentrations ($> 0.3 \text{ mg m}^{-3}$), Model P conforms to the OC4 algorithm (Model L), and for concentrations between > 0.25 and $\leq 0.3 \text{ mg m}^{-3}$ a mixture of the colour index (CI) and the OC4 algorithm (Model L) is used, allowing a smooth transition from the CI to the OC4 with increasing chlorophyll.

3.3. Diffuse attenuation models (K_d)

Algorithms for computing the diffuse attenuation coefficient at 489 nm ($K_d(489)$) are described in the following section. For semi-analytical Models A to J, $K_d(489)$ was computed following Lee et al. (2005), with $a(489)$ and $b_b(489)$

413 computed according to the particular model (A-J) and the solar sun-zenith angle
 414 (θ) as input, such that:

$$K_d(489) = [1 + (0.005\theta)]a(489) + 4.18\{1 - 0.52 \exp[-10.8a(489)]\}b_b(489). \quad (6)$$

415 For semi-analytical Model K, $K_d(489)$ is an output, tied to scattering and total
 416 absorption. In addition to $K_d(489)$ estimates from semi-analytical models, an
 417 empirical algorithm was also incorporated into the comparison (Model Q).

418 3.3.1. Model Q

419 Model Q refers to the NASA empirical algorithm for deriving $K_d(489)$
 420 from SeaWiFS (KD2S). This is a polynomial algorithm that relates the log-
 421 transformed ratio of remote-sensing reflectances (X) to $K_d(489)$. The algorithm
 422 uses a two-band blue-green reflectance ratio to compute X (see Eq. 5), and
 423 $K_d(489)$ is computed following:

$$K_d(489) = 10^{(a_0 + a_1 X + a_2 X^2 + a_3 X^3 + a_4 X^4)} + 0.0166, \quad (7)$$

424 where, $a_0 = -0.8515$, $a_1 = -1.8263$, $a_2 = 1.8714$, $a_3 = -2.4414$ and $a_4 =$
 425 -1.0690 (NASA, 2009).

426 4. Methods

427 4.1. Statistical Tests

428 To test the performance of the in-water algorithms the following univariate
 429 statistical tests were adopted that are commonly used in comparisons between

430 modelled and *in situ* data (e.g. Doney et al., 2009; Friedrichs et al., 2009).

431 4.1.1. Pearson correlation coefficient (r)

432 The correlation coefficient r (also called Pearson's product moment correla-
433 tion) is calculated according to

$$r = \frac{1}{N-1} \sum_{i=1}^N \left[\frac{X_i^M - \left(\frac{1}{N} \sum_{j=1}^N X_j^M \right)}{\left\{ \frac{1}{N-1} \sum_{k=1}^N \left[X_k^M - \left(\frac{1}{N} \sum_{l=1}^N X_l^M \right) \right]^2 \right\}^{1/2}} \right] \left[\frac{X_i^E - \left(\frac{1}{N} \sum_{m=1}^N X_m^E \right)}{\left\{ \frac{1}{N-1} \sum_{n=1}^N \left[X_n^E - \left(\frac{1}{N} \sum_{o=1}^N X_o^E \right) \right]^2 \right\}^{1/2}} \right] \quad (8)$$

434 where, X is the variable and N is the number of samples. The superscript E de-
435 notes the estimated variable (from the model) and the superscript M denotes the
436 measured variable (from NOMAD). Note that the Pearson correlation coefficient
437 assumes a linear relationship between variables and normal distributions. The
438 correlation coefficient may take any value between -1.0 and 1.0.

439 4.1.2. Root Mean Square Error (Ψ)

440 The absolute Root Mean Square Error (Ψ) is calculated according to

$$\Psi = \left[\frac{1}{N} \sum_{i=1}^N \left(X_i^E - X_i^M \right)^2 \right]^{1/2}. \quad (9)$$

441 4.1.3. The bias (δ)

442 The bias between model and measurement can be expressed according to

$$\delta = \frac{1}{N} \sum_{i=1}^N \left(X_i^E - X_i^M \right). \quad (10)$$

443 4.1.4. The centre-pattern Root Mean Square Error (Δ)

444 The absolute centre-pattern (or unbiased) Root Mean Square Error (Δ) is
445 calculated according to

$$\Delta = \left(\frac{1}{N} \sum_{i=1}^N \left\{ \left[X_i^E - \left(\frac{1}{N} \sum_{j=1}^N X_j^E \right) \right] - \left[X_i^M - \left(\frac{1}{N} \sum_{k=1}^N X_k^M \right) \right] \right\}^2 \right)^{1/2}. \quad (11)$$

446 It describes the error of the estimated values with respect to the measured ones,
447 regardless of the average bias between the two distributions. It is related to Ψ
448 and δ according to $\Delta^2 = \Psi^2 - \delta^2$.

449 4.1.5. Slope (S) and Intercept (I) of a Type-2 regression

450 The performance of a model with respect to *in situ* data can be tested us-
451 ing linear regression between the estimated variable (from the model) and the
452 measured variable (*in situ* data), such that

$$X^E = X^M S + I. \quad (12)$$

453 A slope (S) close to one and an intercept (I) close to zero is an indication that the
454 model compares well with the *in situ* data. Type-1 regression typically assumes
455 the dependent variable (*in situ* data) is known infinitely well, when in reality the
456 *in situ* data are also affected by uncertainties (e.g. problems with *in situ* data
457 sampling techniques) that are difficult to quantify. Therefore, we adopted Type-
458 2 regression (Glover et al., 2011, MATLAB function lsqfitma.m), which instead
459 of minimising the vertical distance between independent data and linear fit (as in

460 Type-1 regression), minimises the perpendicular distance between independent
461 data and linear fit.

462 4.1.6. *Percentage of possible retrievals (η)*

463 Considering that algorithms chosen for climate studies should perform rou-
464 tinely, and globally, and should not be a source of more gaps in the data than
465 would be the case if other algorithms were used, the percentage of possible re-
466 trievals (η) is an important criterion that should be considered in the comparison,
467 calculated according to

$$\eta = \frac{N^E}{N^M} 100, \quad (13)$$

468 where N^E represents the number of retrievals using the model and N^M represents
469 the number of *in situ* data points.

470 All statistical tests described above were performed in \log_{10} space, consider-
471 ing the majority of variables are approximately log-normally distributed, with the
472 exception of S_{dg} , γ and $a_{ph}(555)/a_{ph}(443)$ for which the analysis was performed
473 in linear space.

474 4.2. *Quantitative statistical methodology*

475 As with the OC-CCI comparison of atmospheric correction algorithms
476 (Müller et al., Submitted), a points scoring classification was used in the in-water
477 comparison to rank objectively the performance of the algorithms. Each variable
478 was tested independently in the points scoring classification. For each variable,
479 $R_{rs}(\lambda)$ values in the database were used as input to the algorithm to estimate the

variable, the estimated variable was then compared with the corresponding *in situ* value using each statistical test and a score was assigned for each test ranging from zero to two. These tests are described in the following sections. If the algorithm was not capable of estimating the variable, it was given zero points for that test.

In addition, a chi-square test was also performed separately on a selection of the semi-analytical models. This information was used to evaluate the goodness of fit of the computed spectral R_{rs} values compared with the observed values. The samples were only compared when the measured and estimated variables conformed to the following requirements, which represent extreme upper and lower boundaries fixed to avoid the influence of spurious results on the statistical tests (note that algorithms were penalised (Eq. 13) for a higher number of spurious results):

- $C > 0.001$ and $< 200 \text{ mg m}^{-3}$;
- $K_d > a_w$ (Pope and Fry, 1997) and $< 10.0 \text{ m}^{-1}$;
- $a > a_w$ (Pope and Fry, 1997) and $< 10.0 \text{ m}^{-1}$;
- $a_{dg} > 0.0001$ and $< 10.0 \text{ m}^{-1}$;
- $a_{ph} > 0.0001$ and $< 10.0 \text{ m}^{-1}$;
- $b_b > b_{bw}$ (Zhang et al., 2009) and $< 10.0 \text{ m}^{-1}$;
- $\gamma > 0$ and < 4.32 (slope of pure water from Morel, 1974);
- $S_{dg} > 0$ and $< 0.05 \text{ nm}^{-1}$;

- $a_{ph}(555)/a_{ph}(443) > 0$ and < 5.0

The lower boundaries for a_{dg} and a_{ph} were chosen based on the raw uncertainty of a WET-Labs ac9 in waters with low attenuation (WET-Labs, 2012), and lower boundaries for C were based on the absolute accuracy for HPLC detection if all protocols are strictly followed (Aiken et al., 2009). The exclusion of spurious results was conducted on a variable-by-variable basis. For instance, for a given R_{rs} spectra, if a semi-analytical model has one variable (e.g. $a_{ph}(443)$) that falls outside selected boundaries but another (e.g. $a(443)$) that falls within selected boundaries, the former would be excluded and the latter included.

4.2.1. Pearson correlation coefficient (r) test

The r test involved determining whether the r -value for each model was statistically higher or lower than the mean r -value for all models. This was determined using the z_{score} . The z_{score} may be used to determine if two correlation coefficients are statistically different from one another (Cohen and Cohen, 1983). Knowing the r -value for two respective models (say r_1 and r_2 , for model 1 and 2 respectively) and knowing the number of samples used to determine the r -values (say n_1 and n_2) one can determine the z_{score} using the Fisher's r -to- z transformation. Making use of the sample size employed to obtain each coefficient, these z -scores of each r -value (z_1 and z_2) can be used to compute the overall z_{score} (Cohen and Cohen, 1983), such that:

$$z_1 = 0.5 \log\left(\frac{1 + r_1}{1 - r_1}\right), \quad (14)$$

$$z_2 = 0.51 \log\left(\frac{1 + r_2}{1 - r_2}\right), \quad (15)$$

$$z_{score} = \frac{z_1 - z_2}{\{[1/(n_1 - 3)] + [1/(n_2 - 3)]\}^{1/2}}. \quad (16)$$

Having determined the z_{score} , this can be converted into a p -value assuming normal distribution. For the in-water comparison, a two-tailed test was used and if the p -value was <0.05 , the r -values were deemed to be statistically different.

The mean r -value for all models was first determined by averaging the r -value of all the models being tested. The mean number of samples used to compute the r -value, was also determined by averaging all models being tested. The r -value and number of samples of a particular model were then compared with the mean value for all models, so as to determine if the model's r -value was statistically lower, similar or higher than the average value for all models. The following points for each model were awarded accordingly:

- 0 points = r -value for the model tested was statistically lower than the mean r -value for all models.
- 1 point = r -value for the model tested was statistically similar to the mean r -value for all models.
- 2 points = r -value for the model tested was statistically higher than the mean r -value for all models.

539 4.2.2. *Root Mean Square Error (Ψ) and centre-pattern Root Mean Square Error*
540 *(Δ) tests*

541 In addition to computing Ψ and Δ for each model, it is possible to determine
542 the 95% confidence levels in the Ψ and Δ , which provide an indication of how
543 confident one is in Ψ and Δ estimates. The 95% confidence levels can be com-
544 puted from the standard error of the mean percentage and the t -distribution of the
545 sample size. Confidence levels provide a very powerful way of showing differ-
546 ences and similarities between models. If the 95% confidence intervals of two or
547 more models overlap, then it can be assumed that the models have a statistically
548 similar Ψ or Δ .

549 For each model, the Ψ and Δ were computed in addition to their 95% con-
550 fidence intervals. Furthermore, the average Ψ and Δ value for all models tested
551 and the average 95% confidence interval on these values were also calculated.
552 The following points for each model were awarded accordingly:

- 553 • 0 points = Ψ or Δ for the model tested was statistically higher than the
554 mean Ψ or Δ for all models (95% confidence levels did not overlap).
- 555 • 1 point = Ψ or Δ for the model tested was statistically similar to the mean
556 Ψ or Δ for all models (95% confidence levels overlap with mean values).
- 557 • 2 points = Ψ or Δ for the model tested was statistically lower than the mean
558 Ψ or Δ for all models (95% confidence levels did not overlap).

559 Figure 2 shows an example of the points classification for models in the chloro-
560 phyll (C) comparison using Ψ .

561 4.2.3. Bias (δ) test

562 The closer the model bias (δ) is to the reference value of zero implies that
563 the model corresponds well with the *in situ* data. However, a model could have
564 a δ close to the reference value of zero, when compared with another model, but
565 have a much larger 95% confidence interval, implying lower confidence in the
566 retrieved δ . Therefore, the following points classification was introduced for the
567 bias:

- 568 • 0 points = the 95% confidence interval of δ for a particular model is higher
569 than the mean 95% confidence interval for all models. In addition to this,
570 the bias \pm its 95% confidence interval did not overlap with zero \pm the mean
571 95% confidence interval for all models.
- 572 • 1 point = either, the 95% confidence interval of δ for a particular model is
573 lower than the mean 95% confidence interval for all models, or, the bias \pm
574 its 95% confidence interval overlaps with zero \pm the mean 95% confidence
575 interval, but not both cases.
- 576 • 2 points = the 95% confidence interval of δ for a particular model is lower
577 than the mean 95% confidence interval for all models, and, the bias \pm its
578 95% confidence interval overlaps with zero \pm the mean 95% confidence
579 interval.

580 4.2.4. Slope (S) and Intercept (I) test

581 In addition to computing the intercept (I) and the slope (S) from Type-2 re-
582 gression, it is possible to compute the standard deviation on I and S (Glover

et al., 2011, MATLAB function lsqfitma.m). The closer the intercept (I) is to the reference value of zero and the closer the slope (S) is to the reference value of one, the better the fit between variables. However, a model could have an intercept closer to the reference value of zero and a slope closer to the reference value of one, when compared with another model, but have a much larger standard deviation on its retrieved parameters, implying lower confidence in the fit. Therefore, to account for both these possibilities the following points classification was introduced for the slope (S) parameter:

- 0 points = the standard deviation of the S parameter for a particular model is higher than the mean standard deviation for all models. In addition to this, the S parameter \pm its standard deviation does not overlap with one \pm twice the mean standard deviation for all models.
- 1 point = either, the standard deviation of the S parameter for a particular model is lower than the mean standard deviation for all models, or, the S parameter \pm its standard deviation overlaps with one \pm twice the mean standard deviation for all models, but not both cases.
- 2 points = the standard deviation of the S parameter for a particular model is lower than the mean standard deviation for all models, and, the S parameter \pm its standard deviation overlaps with one \pm twice the mean standard deviation for all models.

The following points classification was introduced for intercept (I) parameter:

- 605 • 0 points = the standard deviation of the I parameter for a particular model
 606 is higher than the mean standard deviation for all models. In addition to
 607 this, the I parameter \pm its standard deviation does not overlap with zero \pm
 608 twice the mean standard deviation for all models.
- 609 • 1 point = either, the standard deviation of the I parameter for a particular
 610 model is lower than the mean standard deviation for all models, or, the I
 611 parameter \pm its standard deviation overlaps with zero \pm twice the mean
 612 standard deviation for all models, but not both cases.
- 613 • 2 points = the standard deviation of the I parameter for a particular model
 614 is lower than the mean standard deviation for all models, and, the I param-
 615 eter \pm its standard deviation overlaps with zero \pm twice the mean standard
 616 deviation for all models.

617 4.2.5. *Percentage of possible retrievals (η) test*

618 To compare the percentage of possible retrievals (η) between models, the
 619 average percentage of retrievals for all models was computed in addition to its
 620 standard deviation. The following points criteria were set-up:

- 621 • 0 points = η of a model is less than the mean η of all models – its standard
 622 deviation.
- 623 • 1 point = η of a model overlaps with the mean η for all models \pm its stan-
 624 dard deviation.
- 625 • 2 points = η of a model is greater than the mean η of all models + its
 626 standard deviation.

627 4.2.6. *Total points*

628 To rank the performance of each model with reference to a particular variable,
629 all points were summed over each statistical test. The total score for each model
630 was then normalised by the average score of all models being tested. A score of
631 one indicates the performance of a model is average with respect to all models,
632 a score greater than one indicates a model is performing better than the average
633 and a score less than one indicates the model is performing worse than average.
634 Figure 3 shows a flow-chart illustrating the methodology of the scoring system
635 used to intercompare models. Note that a doubling of points (say from 1 to 2)
636 does not imply an algorithm is twice as good; instead it implies that the difference
637 between the two models is statistically significant.

638 The stability of the scoring system, and the sensitivity of the scores, was
639 tested using the method of bootstrapping (Efron, 1979; Efron and Tibshirani,
640 1993). This involved using sampling with replacement to randomly re-sample
641 the *in situ* data (1000 times) creating 1000 new datasets the same size as the
642 original dataset but not identical. The quantitative statistical methodology was
643 then re-run for each new dataset (Monte-Carlo approach) and from the resulting
644 distribution of scores, a mean score for each model was computed. Additionally,
645 a 2.5% and a 97.5% interval on the bootstrap distribution was taken and assumed
646 to be the error-bars or confidence limits on the mean score for each model, rather
647 than standard deviations on the bootstrap distribution, to avoid misinterpretation
648 of results should the bootstrap distribution not follow a normal distribution or be
649 skewed, for instance from the presence of outliers in the data.

650 4.2.7. Chi-square test

651 In addition to the tests described above, a chi-square (χ^2) test was also used to
652 compare performance of a selection of semi-analytical models. For each semi-
653 analytical model tested, a reconstructed reflectance spectrum was produced in
654 forward mode and compared with the *in situ* reflectance data. This was con-
655 ducted on 1713 samples ($K_d(489)$ database) representative of a broad range of
656 oceanic environments inclusive of the major ocean basins (see Fig. 1). The test
657 is designed to examine how well each semi-analytical model performed at re-
658 producing the observations. The results from this test are not incorporated into
659 the points classification, as some semi-analytical models in the comparison are
660 algebraic (e.g. Models A, B, D and E) thus their χ^2 values equal zero. However,
661 the information is useful to evaluate the performance of those semi-analytical
662 algorithms that are not algebraic (Models C, F, G, H, I, J and K). The chi-square
663 was computed for each of the 1713 spectra using the following formula:

$$\chi^2 = \sum_{i=1}^{N_\lambda} \left[R_{rs}^M(i) - R_{rs}^E(i) \right]^2, \quad (17)$$

664 where, the super-script M is the measured reflectance data, and the super-script
665 E is the estimated reflectance data from the model. The lower the χ^2 is, the better
666 the model reproduces the observed reflectance data.

5. Results

5.1. Chlorophyll comparison

Figure 4 shows results of the quantitative comparison on chlorophyll concentration. What is clear from the scatter plots in Fig. 4 is that all the algorithms perform reasonably at estimating chlorophyll when compared with the *in situ* data ($r > 0.75$). Secondly, a visual qualitative comparison of the scatter plots and the results from the points classification score (bar chart in Fig. 4) reveals that the objective points classification appears to be working consistently, such that the models showing larger discrepancies between modelled and *in situ* data in the scatter plots (e.g. Models C and K) have a low score, and models showing a tighter relationship between modelled and *in situ* data in the scatter plots (e.g. Models L to P) have a higher score.

Results from the classification in Fig. 4 (bar chart) highlight that the empirical chlorophyll models have the highest score (e.g. Model L, M, N and P). This is not surprising considering that many of the *in situ* data used to parameterise these empirical models are not independent of the *in situ* data used here to test these models (see Table 3 and Section 6.1.1 for a discussion of this aspect). However, it is worth noting that Model O, which is the same mathematical equation as Model L, was parameterised using a theoretical model of ocean colour (Morel and Morel, 2001) tuned using data gathered by the Laboratoire d’Océanographie de Villefranche on K_d and chlorophyll (see Morel and Antoine, 2011, for details), data that are independent of the chlorophyll and R_{rs} data used in this comparison. The high score by Model O support the results from Models L, M, and N, in

690 that the empirical (blue-green band-ratio) chlorophyll algorithms perform with
691 a high score in the quantitative comparison. The performance of the empirical
692 algorithms may reflect their immunity to scale errors in R_{rs} data (e.g. band-ratio,
693 see Fig. 14) or errors induced by instrument noise (e.g. band-difference, see Hu
694 et al., 2012).

695 With regard to chlorophyll derived by the semi-analytical algorithms, Mod-
696 els A, G, H and I have a higher score when compared with Models B, D, E and
697 F. However, overlapping error bars from the bootstrap ensemble run, particularly
698 with regard to Model D and E, clearly indicate the difficulty in ranking the per-
699 formance of many of these semi-analytical models objectively. For Models A,
700 G, H and I, error bars from the bootstrap ensemble overlap with the empirical
701 models, suggesting that the performance of these semi-analytical algorithms are
702 comparable with the empirical algorithms in certain conditions. Models C and
703 K perform with low scores, indicating that these semi-analytical models perform
704 less accurately at deriving chlorophyll when compared with the other models in
705 the comparison (Fig. 4).

706 5.2. $K_d(489)$ comparison

707 Figure 5 shows results of the quantitative comparison on $K_d(489)$. All models
708 are seen to capture a high amount of the variability in the $K_d(489)$ *in situ* data (r
709 >0.93). The bar chart indicates empirical Model Q performs with a high points
710 score in the $K_d(489)$ comparison, followed by semi-analytical Models D and E.
711 Models F, I, J and K are shown to perform similarly (slightly above average with
712 scores >1), followed by Models G, H and C. Models A and B have low scores.

Model A shows a systematic overestimation in $K_d(489)$. Considering $a(489)$ and $b_b(489)$ are used as inputs to Eq. (7), this overestimation in $K_d(489)$ associated with Model A can be linked to an overestimation in $b_b(489)$ for this model (see Figure 10) as opposed to the influence of $a(489)$ (see Figure 7).

5.3. The total absorption coefficient ($a(\lambda)$) comparison

Figures 6 and 7 show results of the quantitative intercomparison on $a(\lambda)$. Assessing the scatter plots (Fig. 7), all models capture a high amount of the variability in the *in situ* data at blue and green wavelengths (412-510 nm, $r > 0.87$); at longer wavelengths (e.g. 665 nm), Models A, B, D, and E (all algebraic approaches) have a low score in comparison with the other IOP models in the points classification (Fig. 6). When summing scores over all the wavelengths ($a(\lambda)$ Fig. 6), results from the points classification indicate that, with the exception of Model F which has the highest score, Models C through to K perform with similar scores, as indexed by overlapping error bars. Model A and B have a slightly lower score, which can be attributed to lower scores at longer wavelengths (e.g. Model A and B have a similar score to some models at shorter wavelengths (411, 443 and 489 nm, note the overlapping error bars), but lower scores at longer wavelengths (>510 nm) in Fig. 6). Models A, B, D and E retrieve $a(665)$ directly from $R_{rs}(665)$, consequently when $R_{rs}(665)$ is very low and has a high signal-to-noise ratio (common in oceanic waters), this will result in low quality $a(665)$. However, in such cases, semi-analytical optimisation models (e.g. Models C, F, G, H, I and K) have less dependence on the quality of $R_{rs}(665)$, as $a(665)$ is inferred using a bio-optical model that operates a minimi-

736 sation using wavelengths in blue, green and red regions of the spectrum, often
737 with fixed spectral shapes for the IOPs.

738 5.4. The absorption coefficient of phytoplankton ($a_{ph}(\lambda)$) comparison

739 Figures 6 and 8 show results of the quantitative intercomparison on $a_{ph}(\lambda)$.
740 The results indicate a large range of variability between semi-analytical models.
741 Models A, B, D, and E (algebraic approaches) perform reasonably well at shorter
742 wavelengths (411-489nm), as indexed by a higher points score, but perform less
743 accurately at longer wavelengths (555-665 nm), as indexed by a lower points
744 score. Models C and F through to J alternatively have a higher points score at
745 longer wavelengths (510-665 nm) and lower points score at shorter wavelengths,
746 likely a result of the algebraic approaches performing less accurately at longer
747 wavelengths (555-665 nm). When summing the points across all wavelengths
748 ($a_{ph}(\lambda)$ Fig. 6), Models I and J have the highest scores followed by Models C, G,
749 and H. Model J computes $a_{ph}(\lambda)$ assuming relationships between the chlorophyll
750 concentration of three size-classes of phytoplankton (micro-, nano- and pico-
751 phytoplankton), and their associated specific absorption coefficient (a_{ph}^*), as does
752 Model C during a first iteration to compute b_{bp} and a_{dg} . Models G and H estimate
753 $a_{ph}(\lambda)$ as a linear function of chlorophyll and Model I relates changes in the
754 spectral shape of a_{ph}^* with changes in chlorophyll. Models A and F have an
755 average score (~ 1), in comparison with the other models, with Model K having
756 the lowest score when summing the points across all wavelengths.

757 Figures 6 and 11 show results of the quantitative intercomparison on
758 $a_{ph}(555)/a_{ph}(443)$. Models A and B are seen to perform less accurately at es-

759 estimating $a_{ph}(555)/a_{ph}(443)$, as indexed by a low points score. This can be at-
 760 tributed to the fact that $a_{ph}(555)$ is strongly overestimated by Models A and B
 761 despite performing well at retrieving $a_{ph}(443)$ (Fig. 8), causing an overestimation
 762 of $a_{ph}(555)/a_{ph}(443)$ (Fig. 11). Models C, F, I, and J have the highest scores for
 763 $a_{ph}(555)/a_{ph}(443)$, and it is worth noting that these models tie the spectral shape
 764 of a_{ph} to either the chlorophyll concentration or $a_{ph}(443)$ (Model C only during
 765 a first iteration). Models D, E and K have intermediate scores, as do Models G
 766 and H which assume a fixed spectral shape for a_{ph} (scores of ~ 1). Overlapping
 767 error bars indicate the scores of some of these models are statistically similar.

768 5.5. The absorption coefficient by detrital and dissolved matter ($a_{dg}(\lambda)$) compar- 769 ison

770 Figures 6 and 9 show results of the quantitative intercomparison on $a_{dg}(\lambda)$.
 771 In comparison with $a(\lambda)$ (Fig. 7), the majority of semi-analytical models are
 772 seen to capture a lower amount of the variability in *in situ* $a_{dg}(\lambda)$ ($r \leq 0.88$), in-
 773 dicating lower performance in retrieving $a_{dg}(\lambda)$ in comparison with $a(\lambda)$, at least
 774 for blue and green wavelengths. Slight variations in the performance of the al-
 775 gorithms for each wavelength are observed over the visible spectrum, which is
 776 likely caused by variations in S_{dg} and the spectral shape of a_{ph} between models.
 777 Despite these variations, the points score of all algorithms when summed across
 778 all wavelengths ($a_{dg}(\lambda)$ Fig. 6), is strikingly similar to the performance of the
 779 models at a single wavelength (e.g. $a_{dg}(443)$), highlighting the importance of
 780 correctly estimating the magnitude of a_{dg} at a reference wavelength. However,
 781 it is worth noting that the NOMAD $a_d(\lambda)$ and $a_g(\lambda)$ multi-spectral data were de-

782 developed by fitting an exponential slope to original data on a sample-by-sample
783 basis, to remove moderate noise often resulting from instrument artifacts or poor
784 sample baselines (Werdell, 2005). When summing scores across all wavelengths
785 ($a_{dg}(\lambda)$ Fig. 6), Models D and F have slightly higher scores, followed by Mod-
786 els H, G, E, B, J, A and I. However, with the exception of Models C and K,
787 which have consistently low scores, many models have overlapping error bars
788 indicating statistically similar results.

789 Figures 6 and 11 show results of the quantitative intercomparison on S_{dg} . To
790 compute S_{dg} for each semi-analytical model and *in situ* sample, the spectral a_{dg}
791 results were fitted using an exponential equation between 411-665 nm. What is
792 clear from the scatter plots is that none of the models capture well the variability
793 in S_{dg} ($r < 0.15$, Fig. 11). Models C to F and Model J have a slightly higher
794 score in the points classification when compared with Models A, B, G, H, I and
795 K. The higher points score for Models C to F and J are related to a lower Ψ , Δ
796 and δ for these models (Fig. 11). It is worth noting that Models G, H, and I, have
797 higher S_{dg} (0.018 to 0.0206) than the other models in the comparison.

798 5.6. The total backscattering coefficient ($b_b(\lambda)$) comparison

799 Figures 6 and 10 show results of the quantitative intercomparison on $b_b(\lambda)$.
800 Results indicate that it is difficult to separate the performance of the semi-
801 analytical models at determining $b_b(\lambda)$, as indexed by large error bars on the
802 mean score of the bootstrap distribution. These larger error bars are in part a con-
803 sequence of a lower number of *in situ* samples in the $b_b(\lambda)$ dataset, as compared
804 with the other IOPs. Models A and B display a positive bias (Fig. 10), indicating

805 an overestimation of $b_b(\lambda)$, and Model J appears to underestimate $b_b(\lambda)$ at larger
806 values (Fig. 10). When summing scores across all wavelengths ($b_b(\lambda)$, Fig. 6),
807 Models A, C and K have lower scores and Models D, G, H and J slightly higher
808 scores, when compared with the majority of models.

809 Figures 6 and 11 show results of the quantitative intercomparison on γ . To
810 compute γ for each semi-analytical model, and for the *in situ* data, the spec-
811 tral b_b results were fitted using a power-law equation between 411-665 nm. As
812 with the $b_b(\lambda)$ points classification, it is difficult to separate the performance of
813 some of the algorithms (overlapping error bars). Models D and E have a higher
814 points scores in the γ test (note for these models the slope of b_{bp} was param-
815 eterised using some of the data in NOMAD), followed by Models B, C and F
816 through to J. Models D, E, F, I and J all vary the spectral dependency of par-
817 ticulate backscattering (b_{bp}) as a function of a blue-green ratio, Model J indi-
818 rectly through chlorophyll which is first estimated using a blue-green ratio from
819 Model L. Models G and H assume a constant spectral dependency of particulate
820 backscattering (b_{bp}). Models A and K have a lower score when compared with
821 the other semi-analytical models.

822 5.7. Chi-square tests

823 Figure 12 shows the results from the chi-square (χ^2) test for the non-algebraic
824 semi-analytical models (Models C, F, G, H, I, J, and K). Results indicate that the
825 models with the lowest chi-square are Models I and F, followed by Model K then
826 Models G, H and C. Model J has a higher chi-square when compared to the other
827 models, indicating the agreement between R_{rs} *in situ* and model is lower for this

828 model. For the algorithms that use non-linear optimisation (Models C, F, G, H,
829 I and K) the chi-square results are influenced by both the convergence criteria of
830 the optimisation scheme and the degrees of freedom in the bio-optical model. A
831 more stringent convergence criterion can result in a lower chi-square, but only to
832 an extent that is constrained by the freedom of the model to reproduce observed
833 R_{rs} . The chi-square is also dependent upon the optimisation scheme itself (e.g.
834 Levenberg-Marquardt, Gradient descent, Nelder-Mead method, Quasi-Newton,
835 Trust region), each of which has its advantages and disadvantages (see Mu et al.,
836 2011), how each approach minimises the χ^2 (minimising to the absolute values of
837 R_{rs} , relative values, or even logarithmically transformed values), and the number
838 of wavelengths used in the minimisation.

839 5.8. Overarching comparison of semi-analytical models

840 Figure 13 shows results for the quantitative intercomparison when combining
841 the points score for all variables for each semi-analytical model, then normalising
842 with respect to the mean score. This was conducted in four ways: (i) all points for
843 spectral IOPs ($a(\lambda)$, $a_{dg}(\lambda)$, $a_{ph}(\lambda)$, $b_b(\lambda)$, γ , $a_{ph}(555)/a_{ph}(443)$ and S_{dg}), chloro-
844 phyll (C) and $K_d(489)$; (ii) all points for all spectral IOPs and $K_d(489)$; (iii)
845 all points for all spectral IOPs; (iv) and all points for IOPs from wavelengths
846 411-555 nm. The later was conducted as some algorithms perform poorly at re-
847 trievaling some IOPs at 665 nm (e.g. Model A, B, D, and E) which could have
848 repercussions on the points score for other models (see discussion on this aspect
849 in Section 6.1.2).

850 When combining the scores of all these variables, regardless of approach (i-

iv above), it is evident that Models D to J have higher scores than Models A, B, C and K. It is important to note that despite this, Models A, B, C and K do, in some cases, have higher or comparable scores to Models D to J for particular variables (Fig. 6). Regarding Models D to J, it is very difficult to objectively rank their performance with respect to each other, considering overlapping error bars. Models H and J have a higher points score than Model E in all cases except when summing points for IOPs from wavelengths 411-555 nm. However, in all cases Model E has a statistically similar score to Models D, F, G and I, as indexed by overlapping error bars, and Models F and G have statistically similar scores to Models H and I. Models D to J all have statistically similar scores for IOPs from wavelengths 411-555 nm. Therefore, results from the objective classification indicate that Models D to J perform similarly, when the ensemble of variables are considered. However, as highlighted in Fig 6, the scores of these models vary depending on product and wavelength.

6. Discussion

6.1. Methodological Uncertainties

6.1.1. Data

This paper focuses on the development of a methodology to classify and rank objectively the performance of a variety of in-water bio-optical algorithms. The classification has been applied to a selection of in-water algorithms and the NOMAD *in situ* dataset. We have used the NOMAD dataset as, to our knowledge, it is the most extensive globally-representative dataset of co-located measurements

873 of *in situ* $R_{rs}(\lambda)$ and in-water variables (IOPs, C and $K_d(489)$). To implement
874 the classification requires a large database. Ideally an inter-comparison of this
875 nature should be performed using a database entirely independent of any data
876 used to parameterise the models. In the intercomparison carried out here, it has
877 been difficult to evaluate the impact of the NOMAD dataset on algorithm per-
878 formance, because most algorithms are influenced to some degree by the dataset
879 (see Table 3). The limited availability of *in situ* observations on $R_{rs}(\lambda)$ and in-
880 water variables, coupled with the need for a large database to implement our
881 objective classification has meant that some data used in the comparison are not
882 independent of those used to parameterise many of the models. This was partly
883 addressed using the bootstrap method which allowed for some investigation into
884 the performance of the algorithms in the context of the range of variability in
885 the dataset. However, the work highlights the need for an independent dataset to
886 be developed and used to evaluate algorithms further, to ascertain the extent to
887 which the results are influenced by this issue.

888 Whereas NOMAD is the most extensive global database of *in situ* $R_{rs}(\lambda)$ and
889 in-water variables (IOPs, C and $K_d(489)$), the distribution of measurements in
890 NOMAD is not equivalent to the distribution in the global ocean. Eutrophic wa-
891 ters are over-represented in NOMAD and oligotrophic waters under-represented
892 (Werdell and Bailey, 2005). Ideally, when comparing global bio-optical algo-
893 rithms, a dataset should be used that corresponds approximately to the distribu-
894 tion of measurements in the global ocean, highlighting the need for continued
895 on-going *in situ* campaigns that focus on the areas of the ocean that are under-

896 represented in *in situ* databases, such as the oligotrophic gyres.

897 In the objective classification, the *in situ* datum is essentially deemed to be
898 the truth, whereas, in reality *in situ* data also have associated errors. Measure-
899 ment outliers were minimised using robust quality control procedures adopted
900 in NOMAD (Werdell and Bailey, 2005). However, quantifying these errors is
901 a very difficult task and some variables have a higher level of uncertainty than
902 others. For some of the statistical tests, the measurement errors were partly ac-
903 counted for (e.g. Type-2 regression). Nonetheless, it is recommended that future
904 efforts include uncertainty indices for *in situ* observations.

905 In this study, *in situ* observations of R_{rs} were used as input to the models. It
906 can be assumed that errors in the *in situ* R_{rs} values are small in comparison to
907 satellite-derived R_{rs} . The performance of the algorithms tested may differ when
908 used with data containing higher levels of noise. The tolerance of the bio-optical
909 models to errors in R_{rs} will need to be evaluated further to reflect realistic satellite
910 measurement conditions. This could be done using simulated datasets (e.g. Lee
911 et al., 2010) or satellite and *in situ* match-ups (e.g. Mélin et al., 2005; Bailey
912 and Werdell, 2006; Maritorena et al., 2010). A global database of satellite and *in*
913 *situ* match-ups would also allow for a thorough investigation into the suitability
914 of coupling different in-water bio-optical models with atmospheric correction
915 models. For example, atmospheric-correction models that focus on estimating
916 the spectral-shape of R_{rs} accurately, with low bias, maybe better suited to band-
917 ratio in-water models. Hu et al. (2012) found that band-difference chlorophyll
918 algorithms are less sensitive than band-ratio algorithms to various errors induced

919 by instrument noise and imperfect atmospheric correction in low chlorophyll
920 waters. It is recommended that future efforts investigate potential synergistic
921 benefits of combining different in-water and atmospheric correction models.

922 6.1.2. *Objective classification*

923 The objective classification developed here is a step toward a fully-automated
924 tool for the comparison and development of emerging bio-optical algorithms.
925 The strategy for algorithm selection has to be open to the possibility that better
926 algorithms will emerge in the future, requiring periodic re-evaluations of algo-
927 rithms, adoptions of new algorithms and re-processing of data archives, as and
928 when necessary. The objective classification developed here can aid the quan-
929 titative comparison between emerging and existing algorithms. However, the
930 classification itself may undergo refinement with use and with changing user
931 requirements.

932 There are issues with using the average performance of all models as a base-
933 line from which to compare algorithm performance. If some algorithms perform
934 very poorly this can significantly influence the average performance of all mod-
935 els, to the extent that it becomes difficult to differentiate between the higher per-
936 forming models. This happened for $a(665)$ and $a_{ph}(665)$ (see Fig. 6). Models
937 A, B, D, and E performed poorly, with high Ψ , Δ and δ in comparison with the
938 other models (Fig. 7 and 8) resulting in minimal points for Models A, B, D, and
939 E and maximum points for all other algorithms. Supplementary Fig. S1 shows
940 the $a(665)$ results with and without the inclusion of Models A, B, D and E. When
941 these models are removed from the comparison, it becomes apparent that Model

942 G, H and J have a higher point score than Model C. This issue is to some extent
943 dependent on the number of algorithms being tested. For instance, if one algo-
944 rithm performs poorly it will have a larger effect on the mean of all models when
945 only a small number of algorithms are being compared.

946 It is also important to note that the objective classification was conducted on
947 a variable-by-variable basis. For example, there is no reason why the scores of
948 the individual absorptions (a_{ph} and a_{dg}) should be related to total absorption (a).
949 In Fig 6, Model K has an average score for $a(443)$ but low score for $a_{ph}(443)$ and
950 $a_{dg}(443)$. The performance of Model K impacts the average performance of all
951 models, such that Models G and H have a higher score for $a_{ph}(443)$ and $a_{dg}(443)$
952 than they do for $a(443)$.

953 Another disadvantage of using the average performance of all models as a
954 baseline from which to compare algorithm performance, is that it gives an in-
955 dication only as to the relative performance of each model with respect to the
956 others, and not in absolute terms. For instance, it is clear from the scatter plots
957 (Fig. 5) that $K_d(489)$ is retrieved better by all models than S_{dg} (Fig. 11), yet it is
958 not clear from the scores in the objective classification (Fig. 6). The univariate
959 statistical tests were chosen in the objective classification as they are commonly
960 used in comparisons between modelled and *in situ* data. However, varying the
961 number of statistical tests in the comparison is likely to influence results. Future
962 refinement of the classification may include incorporating additional statistics,
963 or refining the number of statistical tests used, or even weighing the score of the
964 statistics, should one statistic be deemed more important than others.

965 An additional uncertainty is the challenging issue of how to filter the in-
966 fluence of spurious inversion results. Here, we used extreme upper and lower
967 boundaries for each variable to avoid the influence of spurious results on the
968 statistical tests, filtering results if they fall outside the boundaries. For some
969 optimisation models, inversion results are constrained by positive boundaries
970 which differ among approaches and with those used here to filter results. When
971 the boundaries are hit should we consider the results valid or invalid? One may
972 argue that such results are not valid as they are likely to change if the bound-
973 aries assigned by the optimisation scheme change. Setting the boundaries to
974 the same values for all optimisation models, consistent with those used to filter
975 results from other models, could minimise some differences. However, these
976 boundaries are often chosen according to range of data used for parameterisa-
977 tion, which vary among models. There appears to be some subjectivity in the
978 selection of a suitable criterion for filtering spurious inversion results, yet the de-
979 cision may have a large influence on the results of the classification. For future
980 model comparisons, it is recommended that significant efforts be focused toward
981 the development of an objective filter for spurious inversion results.

982 The models tested here differ implicitly in their treatment of uncertainties in
983 the measured R_{rs} values. Band-ratio algorithm assume negligible uncertainties
984 in the blue to green ratios of R_{rs} . Optimisation methods, that neglect certain
985 bands (e.g. Model C), are effectively assuming very large uncertainties in these
986 neglected bands. These differences impose some unavoidable limits on the com-
987 parison. As progress is made in the quantification of uncertainty in R_{rs} (e.g.

988 Moore et al., 2009) treatment of uncertainties in the various models should be-
989 come less diverse.

990 To account for methodological uncertainties in the classification, bootstrap-
991 ping was introduced. This Monte-Carlo approach not only provides a simple
992 method to check the stability of the results, but also offers a straightforward
993 way to derive confidence estimates on the resulting classification (Efron, 1979;
994 Efron and Tibshirani, 1993), which is useful when comparing model perfor-
995 mance. However, bootstrapping can be computationally expensive and cannot
996 offer insight beyond the range of data to which it is applied.

997 6.2. *Implications for algorithm performance and development*

998 What is clear from the results of the comparison is that the performance of
999 each model varied depending on the product and wavelength being tested. Based
1000 on the results in Figures 4, 5, 6 and 12, Table 4 highlights the variables in which
1001 each semi-analytical model (A-L) performed well and less well in the classifica-
1002 tion. This information may be of use to algorithm developers and to users who
1003 are potentially interested in a specific property, as it highlights components in
1004 these models that may require improvement.

1005 Aside from the individual performance of the models, there are variables
1006 for which all models perform reasonably well or less well at retrieval. From the
1007 scatter plots (Fig. 4 to 11) in general, it is apparent that most models perform well
1008 at retrieving $K_d(489)$, $a(411-555)$ and $a_{ph}(443)$. Some algorithms also retrieve b_b
1009 reasonably well. Decomposing a into a_{ph} and a_{dg} is a problem with some models.
1010 An increase in performance of a_{ph} often results in a reduction in performance of

1011 a_{dg} and vice-versa (e.g. see Fig. 6 Models A and B, and Models D and E). In
 1012 general, all models struggle to retrieve $a_{dg}(\lambda)$, as seen in a higher dispersion in
 1013 the $a_{dg}(\lambda)$ scatter plots (Fig. 9) compared with other variables, confirming other
 1014 studies (e.g. Mélin et al., 2007). Many of the models also struggle at retrieving
 1015 $a_{ph}(555)/a_{ph}(443)$ and S_{dg} , since they assume fixed values for these variables
 1016 despite clear variability in the *in situ* data (Fig. 11). As previously highlighted,
 1017 some of these *in situ* variables may have a higher level of measurement error than
 1018 others, which is also dependent on the signal-to-noise ratio of the measurement
 1019 at the wavelength of interest.

1020 Algebraic approaches (Models A, B, D and E) struggle to retrieve reason-
 1021 able results for a and a_{ph} at 665 nm. These algebraic approaches derive the ab-
 1022 sorption coefficients at a specific wavelength directly from measured R_{rs} at that
 1023 wavelength. Typically, for most Case-1 global waters $R_{rs}(665)$ approaches zero,
 1024 due to the dominating effect of water absorption at this wavelength. Therefore,
 1025 direct retrievals of a_{ph} at 665 nm, when there is little a_{ph} signal, are particularly
 1026 challenging using these algebraic approaches. This is further complicated by
 1027 additional inelastic processes (e.g. Raman scattering) that become increasingly
 1028 important at longer wavelengths. Alternatively, many of the optimisation ap-
 1029 proaches operate a minimisation with respect to the to absolute magnitude of R_{rs} .
 1030 For most Case-1 global waters, where $R_{rs}(665)$ approaches zero, $R_{rs}(665)$ has
 1031 lower weight in the optimisation than R_{rs} at shorter wavelengths, meaning that
 1032 retrievals, such as $a_{ph}(665)$, are actually inferred primarily from R_{rs} at shorter
 1033 wavelengths. Under phytoplankton bloom conditions or turbid waters, where

1034 there is a higher signal in $R_{rs}(665)$, it is a different story. Under such conditions,
1035 variables such as $a_{ph}(665)$ could be derived from the measured $R_{rs}(665)$ using
1036 the algebraic approaches (possibly by shifting the reference wavelength further
1037 into the red or near-infrared). It is also likely that optimisation approaches, that
1038 operate a minimisation with respect to the absolute magnitude of R_{rs} , will give
1039 more weight to $R_{rs}(665)$ when deriving $a_{ph}(665)$ in bloom conditions, despite not
1040 deriving $a_{ph}(665)$ directly from $R_{rs}(665)$.

1041 In this comparison, models were tested against a suite of IOPs, $K_d(489)$ and
1042 chlorophyll. It is important to note that many of these models are not designed
1043 for retrieving all these variables. The algebraic QAA model is not intended to
1044 derive IOPs at wavelengths longer than the reference wavelength, and many of
1045 the optimisation algorithms are typically designed to retrieve IOPs at specific
1046 wavelengths assuming an underlying bio-optical model. The advantages and
1047 disadvantages of each approach are, to a certain degree, characteristic of model
1048 design, making built-in biases difficult to avoid in this comparison. Nonetheless,
1049 this comparison has demonstrated that all the algorithms compared have certain
1050 desirable features. Further algorithm improvements could be explored by com-
1051 bining the best features of various algorithms. The NASA GIOP framework is an
1052 ideal platform for such algorithm development, offering users freedom to specify
1053 and compare various optimisation approaches and parameterisations. Alterna-
1054 tively, algorithm improvements may also come from looking outside the current
1055 set of approaches (e.g. Morel and Gentili, 2009; Shanmugam, 2011).

1056 When using semi-analytical approaches to estimate IOPs, it is generally as-

1057 sumed that there is a good closure between the Apparent Optical Properties
 1058 (AOPs) (or quasi-Inherent Optical Properties, such as R_{rs}) and the IOPs them-
 1059 selves. Figure 14 shows a comparison between measured R_{rs} and modelled R_{rs}
 1060 for 87 samples in NOMAD with corresponding R_{rs} , a and b_b at wavelengths
 1061 from 411-555 nm. Modelled R_{rs} in Fig. 14 was reconstructed using *in situ* a
 1062 and b_b and the approximation of Gordon et al. (1988). What is clear from Fig.
 1063 14, is that for the 87 samples used there is an imperfect closure between R_{rs} and
 1064 modelled R_{rs} reconstructed from the *in situ* IOPs. Interestingly, there appears
 1065 to be better closure when reconstructing the shape of R_{rs} from the IOPs. The
 1066 reasons for this lack of closure are likely related to (i) uncertainty or errors in the
 1067 *in situ* measurements themselves (both IOP and R_{rs}) and (ii) errors in the model,
 1068 both of which require further investigation and have implications for algorithm
 1069 development.

1070 6.3. Algorithm selection for climate studies

1071 Figure 13 indicates that when combining results from all variables, semi-
 1072 analytical Models D through to J have higher scores than Models A, B, C and
 1073 K. Depending on the combination of variables (Fig. 13), it is difficult to rank the
 1074 performance of these algorithms, as many of the models have overlapping error
 1075 bars. The selection of suitable algorithms for any project depends not only on
 1076 the quantitative performance of these algorithms, but also their suitability for the
 1077 applications envisaged and the user requirements.

1078 Algorithm selection for climate-change studies should take into considera-
 1079 tion also the development of future ocean-colour products. The detection of

1080 phytoplankton functional types is an emerging area of research (Nair et al., 2008)
 1081 particularly relevant in the context of a changing climate. The spectral shape of
 1082 the phytoplankton absorption coefficient provides an indication of the commu-
 1083 nity structure of phytoplankton (Sathyendranath et al., 2001, 2004; Ciotti et al.,
 1084 2002). To estimate the particle size distribution from satellite data requires mea-
 1085 surements of the spectral slope of particle backscattering (Loisel et al., 2006;
 1086 Kostadinov et al., 2009). The exponential slope of the CDOM coefficient can
 1087 potentially provide information on the proportions of humic and fulvic acids, the
 1088 semi-labile and refractory fractions, photo-degradation status, and the relative
 1089 contribution of a_d to a_{dg} . Bio-optical algorithms that do not allow for variations
 1090 in the spectral shape of these IOPs are unsuitable for development of such prod-
 1091 ucts (nor are they designed with such applications in mind). Accurate retrievals
 1092 of the phytoplankton absorption coefficient at 670 nm have the potential to im-
 1093 prove chlorophyll estimates, considering that absorption at this wavelength is
 1094 less affected by absorption from accessory pigments, and allow for estimates of
 1095 the average size of the phytoplankton (Roy et al., 2010). Algorithms that fail to
 1096 detect $a_{ph}(670)$ will be unsuitable for such purposes. Furthermore, algorithms
 1097 that infer $a_{ph}(670)$ from other wavelengths, or from chlorophyll, are not provid-
 1098 ing the independent information required for such purposes.

1099 Algorithms for climate change studies need to be robust in a changing en-
 1100 vironment. For example, if the phytoplankton community structure changes,
 1101 the alteration in community structure should not interfere with the performance
 1102 of the algorithm at retrieving chlorophyll. Empirical relationships that tie one

1103 property to the next need to be minimised in models, since correlations between
1104 elements of the ecosystem may not be stable in a changing climate. Empirical
1105 relationships are based on observations in the past, often pooling data from mul-
1106 tiple years, which may not be a faithful guide to the future state of the ocean.
1107 If empirical relationships are unavoidable, on-going re-calibration is required
1108 to reduce ambiguity in interpretation of results. A theoretical underpinning of
1109 the empirical models should be established to ascertain sensitivity to possible
1110 climate-related scenarios. Algorithms should also be robust against potential
1111 modifications in relationships between optically-significant constituents, mean-
1112 ing that retrievals of the different contributors to ocean colour should ideally
1113 be independent of each another. This would also facilitate seamless merging
1114 of Case-1 and Case-2 algorithms, considering both water-types are vulnerable
1115 to climate-related change. The different ocean-colour products have to be con-
1116 sistent with each other, in the sense that they close the radiation budget with
1117 minimal error. For instance, the empirical nature of Model J was such that when
1118 combining the individual products the radiation budget was not closed with min-
1119 imal error (Fig. 12).

1120 **7. Summary**

1121 An objective classification has been presented designed to rank the quantita-
1122 tive performance of a suite of bio-optical models based on a variety of univariate
1123 statistics. Eleven semi-analytical models, as well as five empirical chlorophyll al-
1124 gorithms and an empirical diffuse attenuation coefficient algorithm, were ranked
1125 for some 29 variables using the NASA NOMAD dataset. Uncertainty in the

1126 ranking, and sensitivity of the objective classification to the test dataset, were
1127 addressed using a bootstrapping (Monte-Carlo) approach. Results from the clas-
1128 sification suggest that algorithm performance varies depending on the product
1129 and wavelength of interest, and that empirical algorithms in general performed
1130 better in the classification than semi-analytical models at retrieving chlorophyll,
1131 either due to their immunity to scale errors or instrument noise in R_{rs} data, or sim-
1132 ply that data used for model parameterisation were not independent of NOMAD.
1133 However, uncertainty in the classification suggest some semi-analytical algo-
1134 rithms performed comparably to the empirical algorithms at retrieving chloro-
1135 phyll. Methodological uncertainties in the approach were discussed, and indicate
1136 the need for an independent *in situ* dataset for testing models, the need for addi-
1137 tional data in undersampled water types, particularly in oligotrophic waters, and
1138 error quantification of *in situ* data. In addition to testing the quantitative perfor-
1139 mance, algorithm selection for climate change studies need also to consider the
1140 suitability of the algorithm for the purpose and the development of future ocean-
1141 colour products. The objective classification developed here has the potential to
1142 be routinely implemented, for testing the performance of emerging ocean-colour
1143 algorithms and aiding their development.

1144 **8. Acknowledgments**

1145 NOMAD data were contributed by participants in the NASA SIMBIOS Pro-
1146 gram (NRA-96-MTPE-04 and NRA-99-OES-09) and by voluntary contributors.
1147 A cruise name accompanies each data record. Cruise details, including contribu-
1148 tors' names, are available online (<http://seabass.gsfc.nasa.gov/seabasscgi/nomad.cgi>)

1149 using the General Search and Cruise Search utilities to facilitate communication,
1150 collaboration, and acknowledgement. NASA should be commended for the de-
1151 velopment of the NOMAD database and for on-going *in situ* activities. This
1152 work is a contribution to the Ocean Colour Climate Change Initiative of the Eu-
1153 ropean Space Agency and was supported by the UK National Centre for Earth
1154 Observation.

1155 **References**

- 1156 Aiken, J., Pradhan, Y., Barlow, R., Lavender, S., Poulton, A., Holligan, P.,
1157 Hardman-Mountford, N. J., 2009. Phytoplankton pigments and functional
1158 types in the Atlantic Ocean: A decadal assessment, 1995-2005. *Deep Sea Re-*
1159 *search I* 56 (15), 899–917.
- 1160 Bailey, S. W., Werdell, P. J., 2006. A multi-sensor approach for the on-orbit
1161 validation of ocean color satellite data products. *Remote Sensing Environment*
1162 102, 12–23.
- 1163 Brewin, R. J. W., Devred, E., Sathyendranath, S., Hardman-Mountford, N. J.,
1164 Lavender, S. J., 2011. Model of phytoplankton absorption based on three size
1165 classes. *Applied Optics* 50 (2), 4535–4549.
- 1166 Bricaud, A., Babin, M., Claustre, H., Ras, J., Tiéche, F., 2010. Light absorp-
1167 tion properties and absorption budget of Southeast Pacific waters. *Journal of*
1168 *Geophysical Research* 115, C08009.
- 1169 Bricaud, A., Babin, M., Morel, A., Claustre, H., 1995. Variability in the chloro-
1170 phyll specific absorption coefficients of natural phytoplankton: Analysis and
1171 parameterization. *Journal of Geophysical Research* 100, 13,321–13,332.

- 1172 Buiteveld, H., Hakvoort, J. H. M., Donze, M., 1994. The optical properties of
1173 pure water. In: Ocean Optics XII, Proc. SPIE, Vol. 2258, edited by: Jaffe, J.
1174 S., SPIE, Bellingham, Washington, 174-183.
- 1175 Ciotti, A. M., Lewis, M. R., Cullen, J. J., 2002. Assessment of the relation-
1176 ships between dominant cell size in natural phytoplankton communities and
1177 the spectral shape of the absorption coefficient. *Limnology and Oceanography*
1178 47 (2), 404–417.
- 1179 Cohen, J., Cohen, P., 1983. Applied multiple regression/correlation analysis for
1180 the behavioral sciences. L. Erlbaum Associates.
- 1181 Devred, E., Sathyendranath, S., Stuart, V., Platt, T., 2011. A three component
1182 classification of phytoplankton absorption spectra: Applications to ocean-
1183 colour data. *Remote Sensing of Environment* 115 (9), 2255–2266.
- 1184 Doerffer, R., Heymann, K., Schiller, H., 2002. Case 2 water algorithm for the
1185 medium resolution imaging spectrometer (MERIS) on ENVISAT. In: Pro-
1186 ceedings of the ENVISAT validation workshop, 9-13 December 2002, ESA
1187 report.
- 1188 Doerffer, R., Schiller, H., 2000. Neural Network for retrieval of concentrations
1189 of water constituents with the possibility of detecting exceptional out of scope
1190 spectra. In: IEEE 2000 International Geoscience and Remote Sensing Sym-
1191 posium, Honolulu, Hawaii USA, p. 714-717.
- 1192 Doerffer, R., Schiller, H., 2006. The MERIS Case 2 water algorithm. *Interna-
1193 tional Journal of Remote Sensing* 28 (3-4), 517–535.
- 1194 Doney, S. C., Lima, I. D., Moore, J. K., Lindsay, K., Behrenfeld, M. J., West-

1195 berry, T. K., Mahowald, N., M., G. D., Takahashi, T., 2009. Skill metrics for
 1196 confronting global upper ocean ecosystem-biogeochemistry models against
 1197 field and remote sensing data. *Journal of Marine Systems* 76, 95–112.

1198 Efron, B., 1979. Bootstrap methods: another look at the jackknife. *Annals of*
 1199 *Statistics* 7, 1–26.

1200 Efron, B., Tibshirani, R. J., 1993. *An Introduction to the Bootstrap*. Chapman
 1201 and Hall, New York.

1202 Franz, B. A., Werdell, P. J., 2010. A Generalized Framework for Modeling of
 1203 Inherent Optical Properties in Ocean Remote Sensing Applications. In: *Ocean*
 1204 *Optics XX*, Anchorage, Alaska, 27th Sept.-1st Oct. 2010.

1205 Friedrichs, M. A. M., Carr, M.-E., Barber, R. T., Scardi, M., Antoine, D., Arm-
 1206 strong, R. A., Asanuma, I., Behrenfeld, M., Buitenhuis, E. T., Chai, F., Chris-
 1207 tian, J. R., Ciotti, A. M., Doney, S. C., Dowell, M., Dunne, J., Gentili, B.,
 1208 Gregg, W. W., Hoepffner, N., Ishizaka, J., Kameda, T., Lima, I., Marra, J.,
 1209 Mélin, F., Moore, J. K., Morel, A., O'Malley, R. T. O., O'Reilly, J. E., Saba,
 1210 V. S., Schmeltz, M., Smyth, T. J., Tjiputraw, J., Waters, K., Westberry, T. K.,
 1211 Winguth, A., 2009. Assessing the uncertainties of model estimates of primary
 1212 productivity in the tropical Pacific Ocean. *Journal of Marine Systems* 76 (1-2),
 1213 113–133.

1214 Garver, S. A., Siegel, D. A., 1997. Inherent optical property inversion of ocean
 1215 color spectra and its biogeochemical interpretation: 1. Time series from the
 1216 Sargasso Sea. *Journal of Geophysical Research* 102, 18,607–18,625.

1217 GCOS, 2011. Systematic observation requirements from satellite-based data

1218 products for climate. Tech. rep., World Meteorological Organisation (WMO),
 1219 7 bis, avenue de la Paix, CH-1211 Geneva 2, Switzerland.

1220 Glover, D. M., Jenkins, W. J., Doney, S. C., 2011. Modeling Methods for Marine
 1221 Science. Cambridge Univeristy Press.

1222 Gordon, H. R., Brown, O. B., Evans, R. H., Brown, J., Smith, R. C., Baker, K. S.,
 1223 Clark, S. K., 1988. A semianalytic radiance model of ocean color. Journal of
 1224 Geophysical Research 93, 10,909–10,924.

1225 Gordon, H. R., Clark, D. K., Brown, J. W., Brown, O. B., Evans, R. H.,
 1226 Broenkow, W. W., 1983. Phytoplankton pigment concentrations in the Mid-
 1227 dle Atlantic Bight: Comparison of ship determinations and CZCS estimates.
 1228 Applied Optics 22, 20–36.

1229 Guo, L., Hunt, B., Santschi, P., 2001. Effect of dissolved organic matter on the
 1230 uptake of trace metals by American oysters. Environmental Science and Tech-
 1231 nology 35 (5), 885–893.

1232 Hu, C., Lee, Z., Franz, B., 2012. Chlorophyll a algorithms for oligotrophic
 1233 oceans: A novel approach based on three-band reflectance difference. Jour-
 1234 nal of Geophysical Research 117, C01011.

1235 Huot, Y., Morel, A., Twardowski, M. S., Stramski, D., Reynolds, R. A., 2008.
 1236 Particle optical backscattering along a chlorophyll gradient in the upper layer
 1237 of the eastern South Pacific Ocean. Biogeosciences 5, 495–507.

1238 IOCCG, 2000. Remote Sensing of Ocean Colour in Coastal, and Other Opti-
 1239 cally Complex waters. Tech. rep., Sathyendranath, S. (e.d.), Reports of the
 1240 International Ocean-Colour Coordinating Group, No. 3, IOCCG, Dartmouth,

1241 Canada.

1242 IOCCG, 2006. Remote Sensing of Inherent Optical Properties: Fundamentals,
 1243 Tests of Algorithms, and Applications. Tech. rep., Lee, Z. P. (e.d.), Reports
 1244 of the International Ocean-Colour Coordinating Group, No. 5, IOCCG, Dart-
 1245 mouth, Canada.

1246 Kostadinov, T. S., Siegel, D. A., Maritorena, S., 2009. Retrieval of the particle
 1247 size distribution from satellite ocean color observations. *Journal of Geophys-
 1248 cal Research* 114, C09015.

1249 Lavender, S., Pinkerton, M. H., Morales, J. F., Aiken, J., Moore, G. F., 2004. Sea-
 1250 WiFS validation in European coastal waters using optical and bio-geochemical
 1251 measurements. *International Journal of Remote Sensing* 25 (7-8), 1481–1488.

1252 Lee, Z., Arnone, R., Hu, C., Werdell, P., Lubac, B., 2010. Uncertainties of op-
 1253 tical parameters and their propagations in an analytical ocean color inversion
 1254 algorithm. *Applied Optics* 49 (3), 369–381.

1255 Lee, Z., Carder, K. L., Arnone, R. A., 2002. Deriving inherent optical properties
 1256 from water color: a multiband quasi-analytical algorithm for optically deep
 1257 waters. *Applied Optics* 41 (27), 5755–5772.

1258 Lee, Z., Carder, K. L., Mobley, C. D., Steward, R. G., Patch, J. S., 1998. Hy-
 1259 perspectral remote sensing for shallow waters. 1. A semianalytical model. *Ap-
 1260 plied Optics* 37 (27), 6329–6338.

1261 Lee, Z., Carder, K. L., Mobley, C. D., Steward, R. G., Patch, J. S., 1999. Hy-
 1262 perspectral remote sensing for shallow waters. 2. Deriving bottom depths and
 1263 water properties by optimization. *Applied Optics* 38, 3831–3843.

1264 Lee, Z., Du, K., Arnone, R., 2005. A model for the diffuse attenuation coefficient
1265 of downwelling irradiance. *Journal of Geophysical Research* 110, C02016.

1266 Lee, Z., Lubac, B., Werdell, P. J., Arnone, R., 2009. An Up-
1267 date of the Quasi-Analytical Algorithm (QAA_v5). Tech. rep., In-
1268 ternational Ocean Colour Coordinating Group (IOCCG) Online:
1269 <http://www.ioccg.org/groups/software.html> (assessed 10/02/2012).

1270 Loisel, H., Nicolas, J.-M., Sciandra, A., Stramski, D., A., P., 2006. Spectral
1271 dependency of optical backscattering by marine particles from satellite remote
1272 sensing of the global ocean. *Journal of Geophysical Research* 111, C09024.

1273 Maritorena, S., Fanton d'Andon, O. H., Mangin, A., Siegel, D. A., 2010. Merged
1274 satellite ocean color data products using a bio-optical model: Characteristics,
1275 benefits and issues. *Remote Sensing Environment* 114, 1791–1804.

1276 Maritorena, S., Siegel, D. A., Peterson, A. R., 2002. Optimization of a semi-
1277 analytical ocean color model for global-scale applications. *Applied Optics*
1278 41 (15), 2705–2714.

1279 Mélin, F., Berthon, J.-F., Zibordi, G., 2005. Assessment of apparent and inher-
1280 ent optical properties derived from SeaWiFS with field data. *Remote Sensing*
1281 *Environment* 97, 540–553.

1282 Mélin, F., Zibordi, G., Berthon, J.-F., 2007. Assessment of satellite ocean color
1283 products at a coastal site. *Remote Sensing Environment* 110, 192–215.

1284 Moore, T. S., Campbell, J. W., Dowell, M. D., 2009. A class-based approach to
1285 characterizing and mapping the uncertainty of the MODIS ocean chlorophyll
1286 product. *Remote Sensing Environment* 113, 2424–2430.

1287 Morel, A., 1974. Optical properties of pure water and pure seawater. In: Jerlov,
1288 N. G., Steemann Nielsen, E. (Eds.), Optical Aspects of Oceanography. Academic,
1289 San Diego, California, pp. 1–24.

1290 Morel, A., 1980. In-water and remote measurements of ocean color. *Boundary*
1291 *Layer Meteorology* 18, 177–201.

1292 Morel, A., 2009. Are the empirical relationships describing the bio-optical prop-
1293 erties of case 1 waters consistent and internally compatible? *Journal of Geo-*
1294 *physical Research* 114, C01016.

1295 Morel, A., Antoine, D., 2011. MERIS Algorithm Theoretical Basis Documents
1296 (ATBD 2.9) - Pigment Index Retrieval in CASE 1 Waters (PO-TN-MEL-GS-
1297 0005), Issue 4, July 2011. Tech. rep., Laboratoire d’Océanographie de Ville-
1298 franche (LOV), MERIS ESL, ACRI-ST.

1299 Morel, A., Antoine, D., Gentili, B., 2002. Bidirectional reflectance of oceanic
1300 waters: accounting for Raman emission and varying particle scattering phase
1301 function. *Applied Optics* 41, 6289–6306.

1302 Morel, A., Gentili, B., 2009. A simple band ratio technique to quantify the col-
1303 ored dissolved and detrital organic material from ocean color remotely sensed
1304 data. *Remote Sensing Environment* 113, 998–1011.

1305 Morel, A., Huot, Y., Gentili, B., Werdell, P. J., Hooker, S. B., Franz, B. A.,
1306 2007. Examining the consistency of products derived from various ocean color
1307 sensors in open ocean (case 1) waters in the perspective of a multi-sensor
1308 approach. *Remote Sensing of Environment* 111, 69–88.

1309 Morel, A., Maritorena, S., 2001. Bio-optical properties of oceanic waters: A

1310 reappraisal. *Journal of Geophysical Research* 106 (C4), 7163–7180.

1311 Morel, A., Prieur, L., 1977. Analysis of variations in ocean color. *Limnology and*
 1312 *Oceanography* 22, 709–722.

1313 Mu, X., Shen, Q and, L. Z.-L., Yan, G., & Sobrino, J. A., 2011. A compari-
 1314 son of different optimization algorithms for retrieving aerosol optical depths
 1315 from satellite data: an example of using a dual-angle algorithm. *International*
 1316 *Journal of Remote Sensing* 32 (24), 8949–8968.

1317 Müller, D., Krasemann, H., Brewin, R. J. W., Brockmann, C., Deschamps, P.-
 1318 Y., Doerffer, R., Fomferra, N., Franz, B. A., Grant, M. G., Groom, S., Mélin,
 1319 F., Platt, T., Regner, P., Sathyendranath, S., Steinmetz, F., Swinton, J., Sub-
 1320 mitted. The Ocean Colour Climate Change Initiative: I. A methodology for
 1321 assessing atmospheric correction processors based on in-situ measurements.
 1322 *Remote Sensing Environment*.

1323 Nair, A., Sathyendranath, S., Platt, T., Morales, J., Stuart, V., Forget, M.-H.,
 1324 Devred, E., Bouman, H., 2008. Remote sensing of phytoplankton functional
 1325 types. *Remote Sensing of Environment* 112 (8), 3366–3375.

1326 NASA, June 2009. Diffuse attenuation coefficient (KD) for downwelling irradi-
 1327 ance at 490-nm.
 1328 URL <http://oceancolor.gsfc.nasa.gov/REPROCESSING/R2009/kdv4/>

1329 NASA, March 2010. Ocean Color Chlorophyll (OC) v6.
 1330 URL <http://oceancolor.gsfc.nasa.gov/REPROCESSING/R2009/ocv6/>

1331 Nelson, N. B., Siegel, D. A., 2013. The Global Distribution and Dynamics of
 1332 Chromophoric Dissolved Organic Matter. *Annual Review of Marine Science*

1333 5, 20.1–20.3.

1334 O'Reilly, J. E., Maritorena, S., Mitchell, B. G., Siegel, D. A., Carder, K. L.,
 1335 Garver, S. A., Kahru, M., McClain, C., 1998. Ocean chlorophyll algorithms
 1336 for SeaWiFS. *Journal of Geophysical Research* 103 (C11), 24,937–24,953.

1337 O'Reilly, J. E., Maritorena, S., Siegel, D. and O'Brien, M. C., Toole,
 1338 D. and Mitchell, B. G., Kahru, M., Chavez, F. P., Strutton, P., Cota, G., Hooker,
 1339 S. B., McClain, C. R., Carder, K. L., Muller-Karger, F., Harding, L., Magnu-
 1340 son, A., Phinney, D., Moore, G. F., Aiken, J., Arrigo, K. R., Letelier, R., Cul-
 1341 ver, M., 2000. Ocean color chlorophyll a algorithms for SeaWiFS, OC2, and
 1342 OC4:. Tech. rep., In: Hooker, S. B., Firestone, E. R. (Eds.), *SeaWiFS Post-*
 1343 *launch Technical Report Series. Vol. 11. SeaWiFS Postlaunch Calibration and*
 1344 *Validation Analyses, Part 3. NASA, Goddard Space Flight Center, Greenbelt,*
 1345 *Maryland. 9-23.*

1346 Pope, R., Fry, E., 1997. Absorption spectrum (380-700 nm) of pure water. II.
 1347 Integrating cavity measurements. *Applied Optics* 36 (33), 8710–8723.

1348 Roy, S., Sathyendranath, S., Platt, T., 2010. Retrieval of phytoplankton size from
 1349 bio-optical measurements: theory and applications. *Journal of the Royal So-*
 1350 *ciety Interface* 8 (58), 650–660.

1351 Santschi, P., Lenhart, J., Honeyman, B., 1997. Heterogeneous processes affecting
 1352 trace contaminant distribution in estuaries: the role of natural organic matter.
 1353 *Marine Chemistry* 58 (1-2), 99–125.

1354 Sathyendranath, S., Platt, T., 1997. Analytic model of ocean color. *Applied Op-*
 1355 *tics* 36, 2620–2629.

1356 Sathyendranath, S., Stuart, V., Cota, G., Maas, H., Platt, T., 2001. Remote sens-
 1357 ing of phytoplankton pigments: a comparison of empirical and theoretical
 1358 approaches. *International Journal of Remote Sensing* 22, 249–273.

1359 Sathyendranath, S., Watts, L., Devred, E., Platt, T., Caverhill, C., Maass, H.,
 1360 2004. Discrimination of diatoms from other phytoplankton using ocean-colour
 1361 data. *Marine Ecological Progress Series* 272, 59–68.

1362 Shanmugam, P., 2011. New models for retrieving and partitioning the colored
 1363 dissolved organic matter in the global ocean: Implications for remote sensing.
 1364 *Remote Sensing Environment* 115, 1501–1521.

1365 Smyth, T. J., Moore, G. F., Hirata, T., Aiken, J., 2006. Semianalytical model for
 1366 the derivation of ocean color inherent optical properties: description, imple-
 1367 mentation, and performance assessment. *Applied Optics* 45 (31), 8116–8131.

1368 Werdell, P. J., 2005. An evaluation of inherent optical property data
 1369 for inclusion in the NASA bio-Optical Marine Algorithm Data
 1370 set,” NASA Ocean Biology Processing Group Science Systems and
 1371 Applications, Inc. Document Version 1.1, corresponding to NO-
 1372 MAD Version 1.3 19th September 2005. Tech. rep., NASA: Online
 1373 http://seabass.gsfc.nasa.gov/seabass/data/werdell_nomad_iop_qc.pdf (2005).

1374 Werdell, P. J., 2009. Global bio-optical algorithms for ocean color satellite ap-
 1375 plications. *Eos Trans. AGU* 90 (1), 4.

1376 Werdell, P. J., Bailey, S. W., 2005. An improved in-situ bio-optical data set for
 1377 ocean colour algorithm development and satellite data production validation.
 1378 *Remote Sensing Environment* 98, 122–140.

1379 Werdell, P. J., Franz, B. A., Bailey, S. W., Feldman, G. C., Boss, E., Brando,
1380 V. E., Dowell, M., Hirata, T., Lavender, S. J., Lee, Z., Loisel, H., Maritorena,
1381 S., Mélin, F., Moore, T. S., Smyth, T. J., Antoine, D., Devred, E., d'Andon, O.
1382 H. F., Mangin, A., 2013. Generalized ocean color inversion model for retriev-
1383 ing marine inherent optical properties. *Applied Optics* 52, 2019–2037.
1384 WET-Labs, 2012. An introduction to in-situ absorption and attenuation meters.
1385 URL <http://www.wetlabs.com/technicalnotes/technoteindex.htm>
1386 Zhang, X., Hu, L., He, M.-X., 2009. Scattering by pure seawater: Effect of
1387 salinity. *Optics Express* 17, 5698–5710.

Table 1: Variables tested in the in-water comparison.

Abbreviation	Variable	Usage	Unit
$L_w(\lambda)$	Spectral water-leaving radiance	Input	$\text{uW cm}^{-2} \text{ nm}^{-1} \text{ sr}^{-1}$
$E_s(\lambda)$	Spectral surface irradiance	Input	$\text{uW cm}^{-2} \text{ nm}^{-1} \text{ sr}^{-1}$
$R_{rs}(\lambda)$	Remote sensing reflectance ($L_w(\lambda)/E_s(\lambda)$)	Input	sr^{-1}
θ	Solar sun-zenith angle	Input [#]	Degrees
C	HPLC chlorophyll-a concentration	Output	mg m^{-3}
$K_d(489)$	Diffuse downwelling irradiance coefficient at 489 nm	Output	m^{-1}
$a(\lambda)$	Total absorption coefficient	Output	m^{-1}
$a_{ph}(\lambda)$	Phytoplankton absorption coefficient	Output	m^{-1}
$a_{dg}(\lambda)$	Dissolved (gelbstoff) and detrital (non-algal) absorption coefficient	Output	m^{-1}
S_{dg}	Exponential slope of a_{dg} with wavelength ^{\$}	Output	nm^{-1}
$a_{ph}(555)/a_{ph}(443)$	Index of spectral shape in a_{ph}	Output	Dimensionless
$b_b(\lambda)$	Total backscattering coefficient	Output	m^{-1}
γ	Power slope of b_b with wavelength &	Output	Dimensionless

^{\$} Computed from fitting an exponential function to *in situ data* and model.

& Computed from fitting a power function to *in situ* data and model.

[#] Solar sun-zenith angle was used as input to some of the semi-analytical models and for estimating $K_d(489)$ (see Eq. (6)).
 λ = wavelength.

Table 2: Model output variables.

Model	Output variable										Reference
	$K_d(489)$	C	$a(\lambda)$	$a_{ph}(\lambda)$	$a_{dg}(\lambda)$	$b_b(\lambda)$	γ	S_{dg}	$a_{ph}(555)/a_{ph}(443)$		
A	\times^*	$\times^{\$}$	\times	\times	\times	\times	\times	\times	\times		Smyth et al. (2006)
B	\times^*	$\times^{\$}$	\times	\times	\times	\times	\times	\times	\times		Smyth et al. (2006)
C	\times^*	\times	\times	\times	\times	\times	\times	\times	\times		Devred et al. (2011)
D	\times^*	$\times^{\$}$	\times	\times	\times	\times	\times	\times	\times		Lee et al. (2002)
E	\times^*	$\times^{\$}$	\times	\times	\times	\times	\times	\times	\times		Lee et al. (2009)
F	\times^*	$\times^{\$}$	\times	\times	\times	\times	\times	\times	\times		Lee et al. (1998, 1999)
G	\times^*	\times	\times	\times	\times	\times	\times	\times	\times		Maritorena et al. (2002)
H	\times^*	\times	\times	\times	\times	\times	\times	\times	\times		Maritorena et al. (2002)
I	\times^*	\times	\times	\times	\times	\times	\times	\times	\times		Werdell et al. (2013)
J	\times^*		\times	\times	\times	\times	\times	\times	\times	see [#]	
K	\times^*	\times	\times	\times	\times	\times	\times	\times	\times	Doerffer and Schiller (2000)	
L		\times								O'Reilly et al. (2000)	
M		\times								O'Reilly et al. (2000)	
N		\times								O'Reilly et al. (2000)	
O		\times								Morel et al. (2007)	
P		\times								Hu et al. (2012)	
Q	\times									NASA (2009)	

* Computed following Eq. (6) with θ , $a(489)$ and $b_b(489)$ as input from the model.

$\$$ Computed following Eq. (1) with $a_{ph}(443)$ as input from the model.

[#] This model represents a Case-1 approach that uses Model L as input. The model computes IOPs as a function of C through combining relationships proposed by: Gordon et al. (1983); Buiteveld et al. (1994); Pope and Fry (1997); Huot et al. (2008); Morel (2009); Bricaud et al. (2010); Brewin et al. (2011).

Table 3: Summary of models used in the comparison.

Model	Approach	Method	Input R_{rs} wavelengths [#]	NOMAD Independence ^{\$}
A	Semi-analytical	Algebraic	411, 443, 489, 510, 555, 665	1
B	Semi-analytical	Algebraic	411, 443, 489, 510, 555, 665	2
C	Semi-analytical	Optimisation	443, 489, 510, 555	2
D	Semi-analytical	Algebraic	411, 443, 489, 510, 555, 665	1
E	Semi-analytical	Algebraic	411, 443, 489, 510, 555, 665	2
F	Semi-analytical	Optimisation	411, 443, 489, 510, 555, 665	1
G	Semi-analytical	Optimisation	411, 443, 489, 510, 555, 665	1
H	Semi-analytical	Optimisation	411, 443, 489, 510, 555, 665	1
I	Semi-analytical	Optimisation	411, 443, 489, 510, 555, 665	1
J	Semi-analytical	Band-ratio	443, 489, 510, 555	2
K	Semi-analytical	Optimisation	411, 443, 489, 510, 555, 665	1
L	Empirical	Band-ratio	443, 489, 510, 555	3
M	Empirical	Band-ratio	443, 489, 555	3
N	Empirical	Band-ratio	489, 555	3
O	Empirical	Band-ratio	443, 489, 510, 555	1
P	Empirical	Band-ratio / CI*	443, 489, 510, 555, 665	3
Q	Empirical	Band-ratio	489, 555	3

[#] Wavelengths used that are available in the comparison.

^{\$} Qualitative assessment of algorithm independence to NOMAD: 1 = NOMAD dataset has a small influence on model parameterisation; 2 = NOMAD dataset has some influence on model parameterisation; 3 = NOMAD dataset has a large influence on model parameterisation.

* CI refers to a colour index defined as the difference between R_{rs} in the green region of the visible spectrum and a reference formed linearly between R_{rs} in the blue and red region of the visible spectrum.

Table 4: Performance of semi-analytical models in the objective classification.

Model	Higher performance	Lower performance
A	$a_{ph}(411-510), b_b(411), C$	$a(\lambda), a_{ph}(555-665), a_{ph}(555)/a_{ph}(443), b_b(510-665), \gamma, K_d(489)$
B	$a_{ph}(411-443), a_{dg}(665)$	$a(489-665), a_{ph}(555-665), a_{ph}(555)/a_{ph}(443), K_d(489)$
C	$a_{ph}(510-555), a_{ph}(555)/a_{ph}(443), S_{dg}$	$a(411-443), a_{dg}(\lambda), b_b(\lambda), C$
D	$a(411-555), a_{dg}(\lambda), b_b(\lambda), S_{dg}, \gamma, K_d(489)$	$a(665), a_{ph}(510-665)$
E	$a(411-555), a_{ph}(411-443), \gamma, K_d(489)$	$a(665), a_{ph}(510-665)$
F	$a(\lambda), a_{dg}(\lambda), S_{dg}, a_{ph}(555)/a_{ph}(443), \chi^2$	$a_{ph}(411), b_b(665)$
G	$a(665), a_{ph}(510-665), b_b(\lambda), C$	$a(443-489)$
H	$a(510-665), a_{ph}(510-665), b_b(\lambda), C$	$a(443)$
I	$a(665), a_{ph}(\lambda), a_{ph}(555)/a_{ph}(443), C, \chi^2$	$a_{dg}(555-665)$
J	$a(665), a_{ph}(\lambda), a_{ph}(555)/a_{ph}(443), S_{dg}$	χ^2
K	$a(665), a_{ph}(555)/a_{ph}(443)$	$a(489-510), a_{ph}(\lambda), a_{dg}(\lambda), b_b(411-510), S_{dg}, \gamma, C$

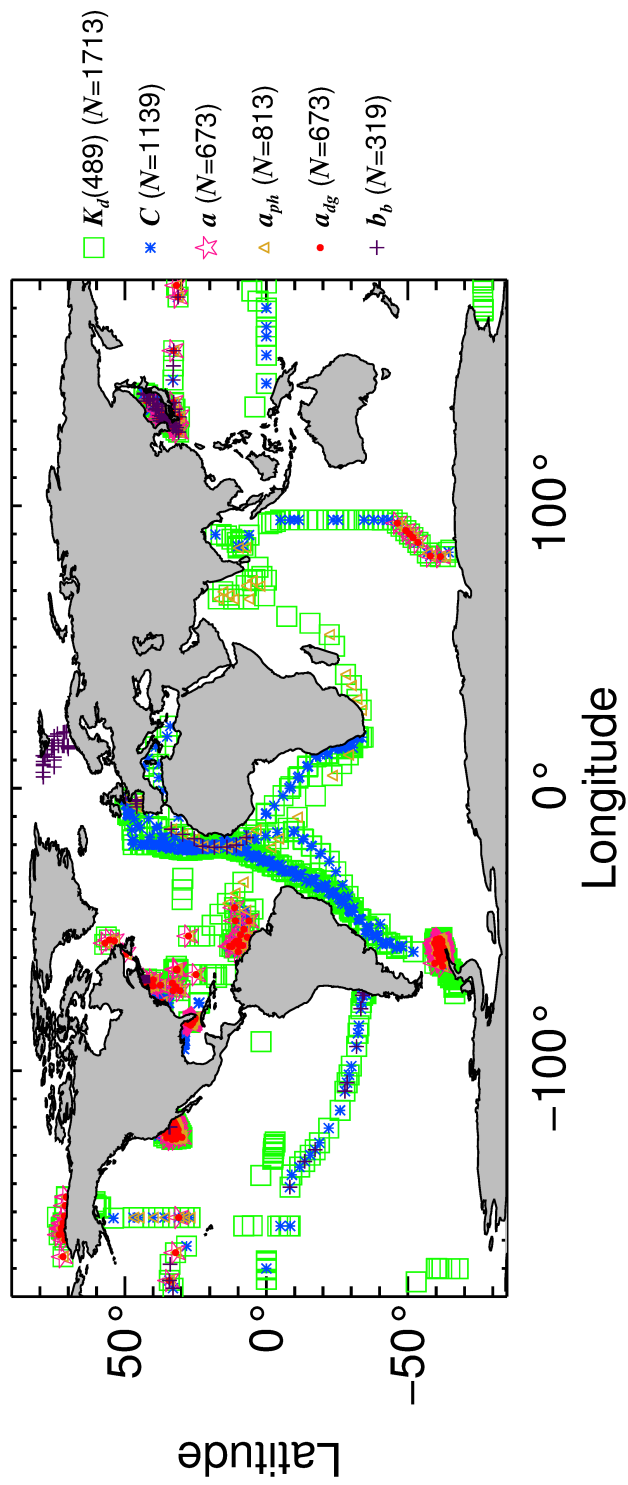


Figure 1: NOMAD *in situ* data used in the study (N = number of samples).

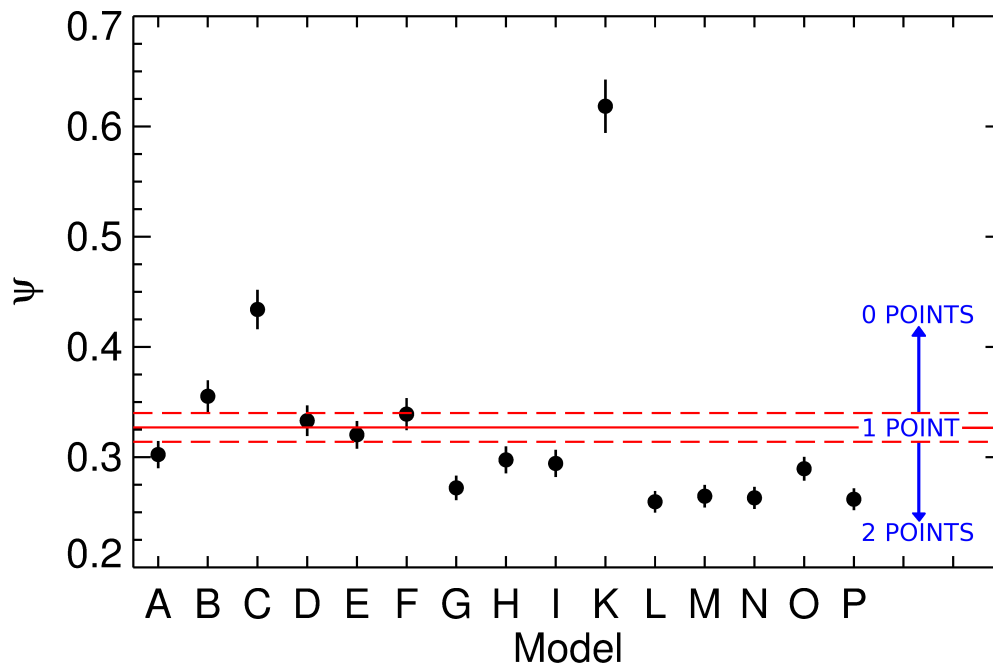


Figure 2: An example of the points classification for a number of models tested in the chlorophyll comparison using the Root Mean Square Error (Ψ). Red solid line represents the mean Ψ for all models and dashed red lines represent the mean $\Psi \pm$ mean 95% confidence intervals. The Ψ of each model is shown by the filled black circle and the black lines represent the Ψ of each model \pm 95% confidence intervals.

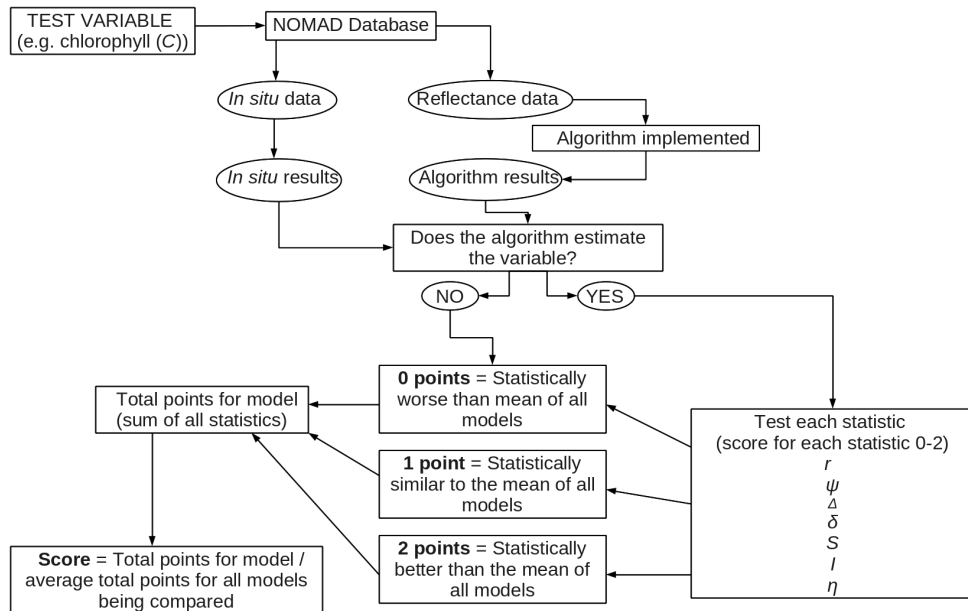


Figure 3: Flow chart illustrating the methodology of the scoring system.

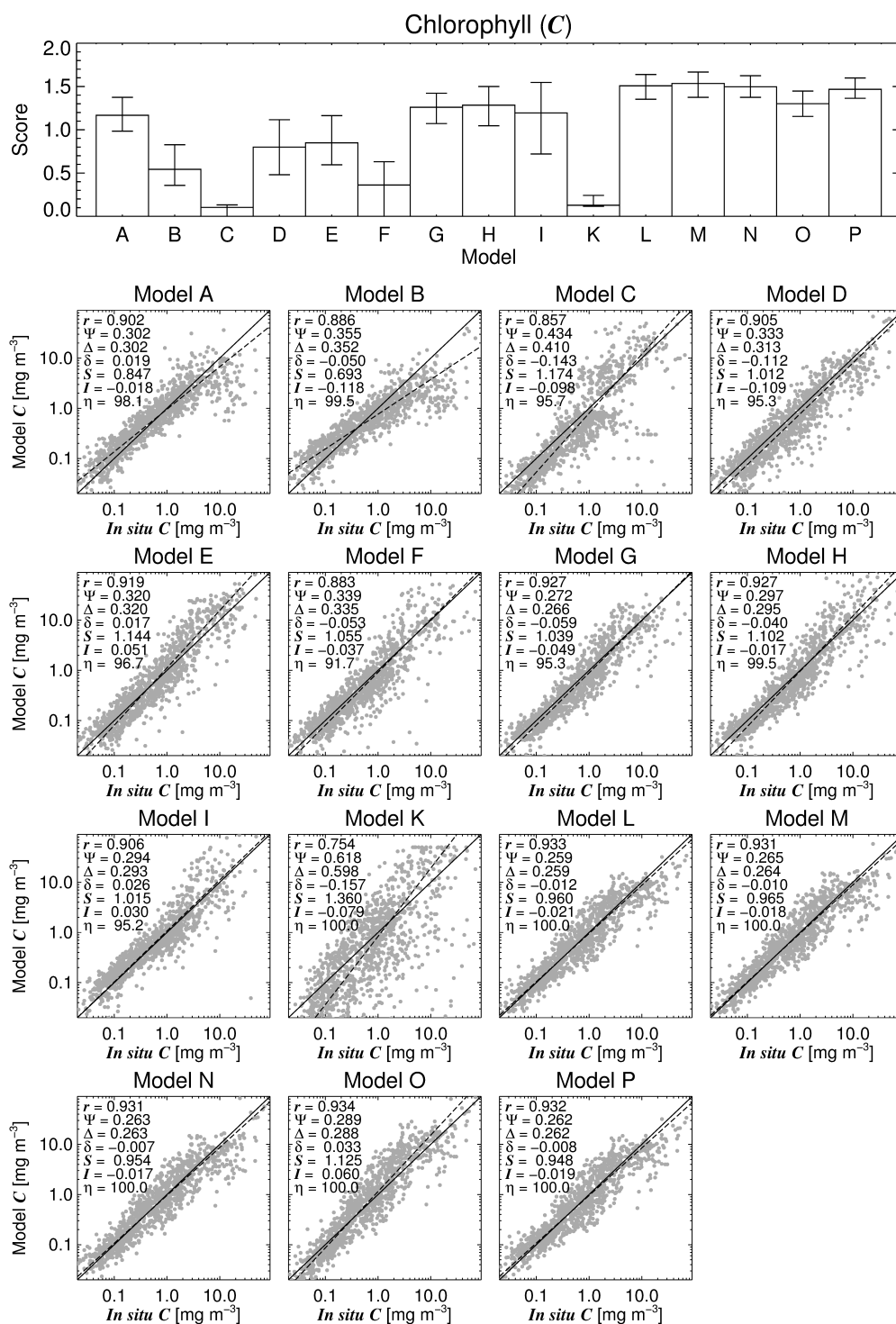


Figure 4: Results from the chlorophyll (C) model comparison.

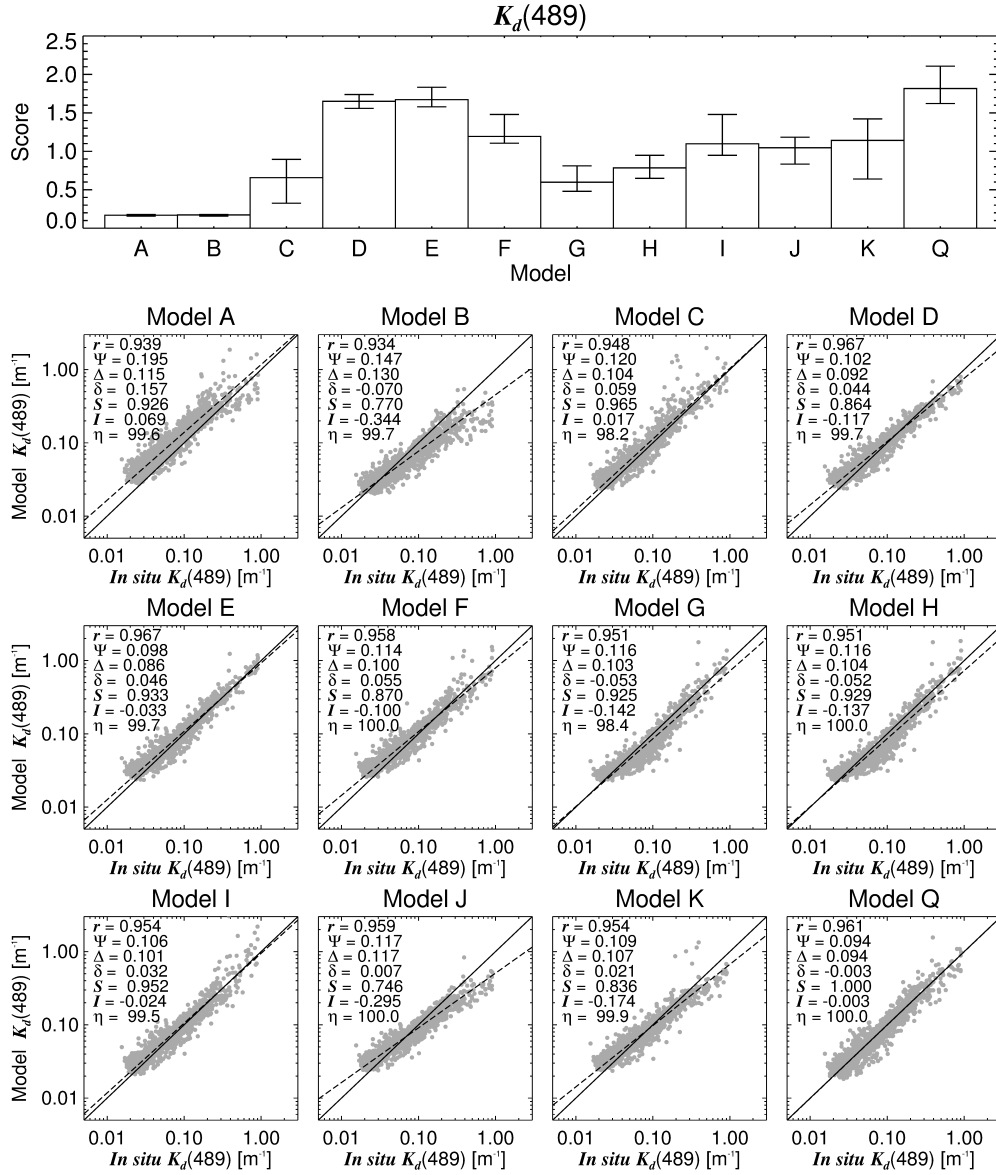


Figure 5: Results from the diffuse attenuation coefficient at 489 nm ($K_d(489)$) model comparison.

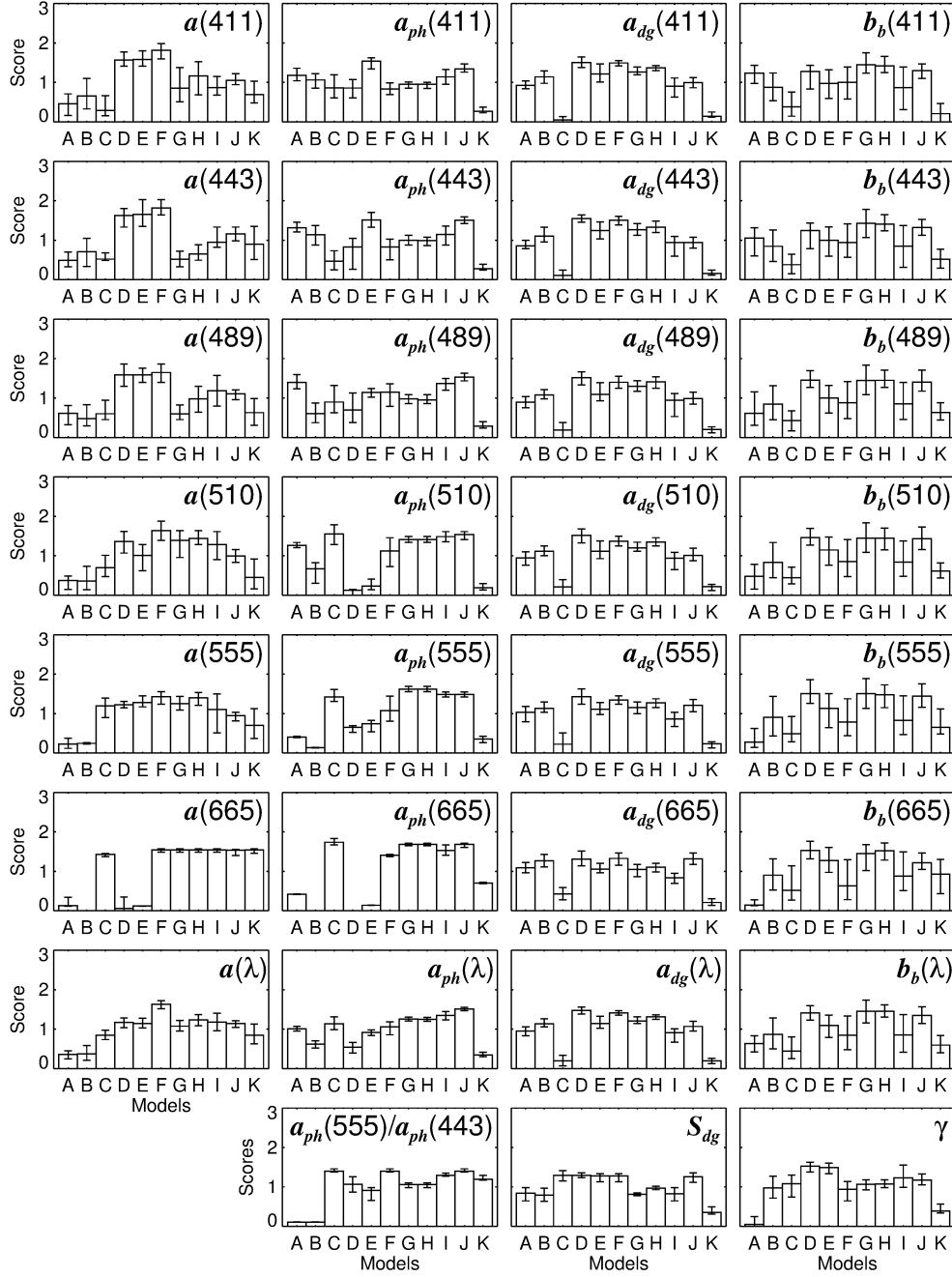


Figure 6: Results of the semi-analytical models at retrieving Inherent Optical Properties (IOP) according to the points classification.

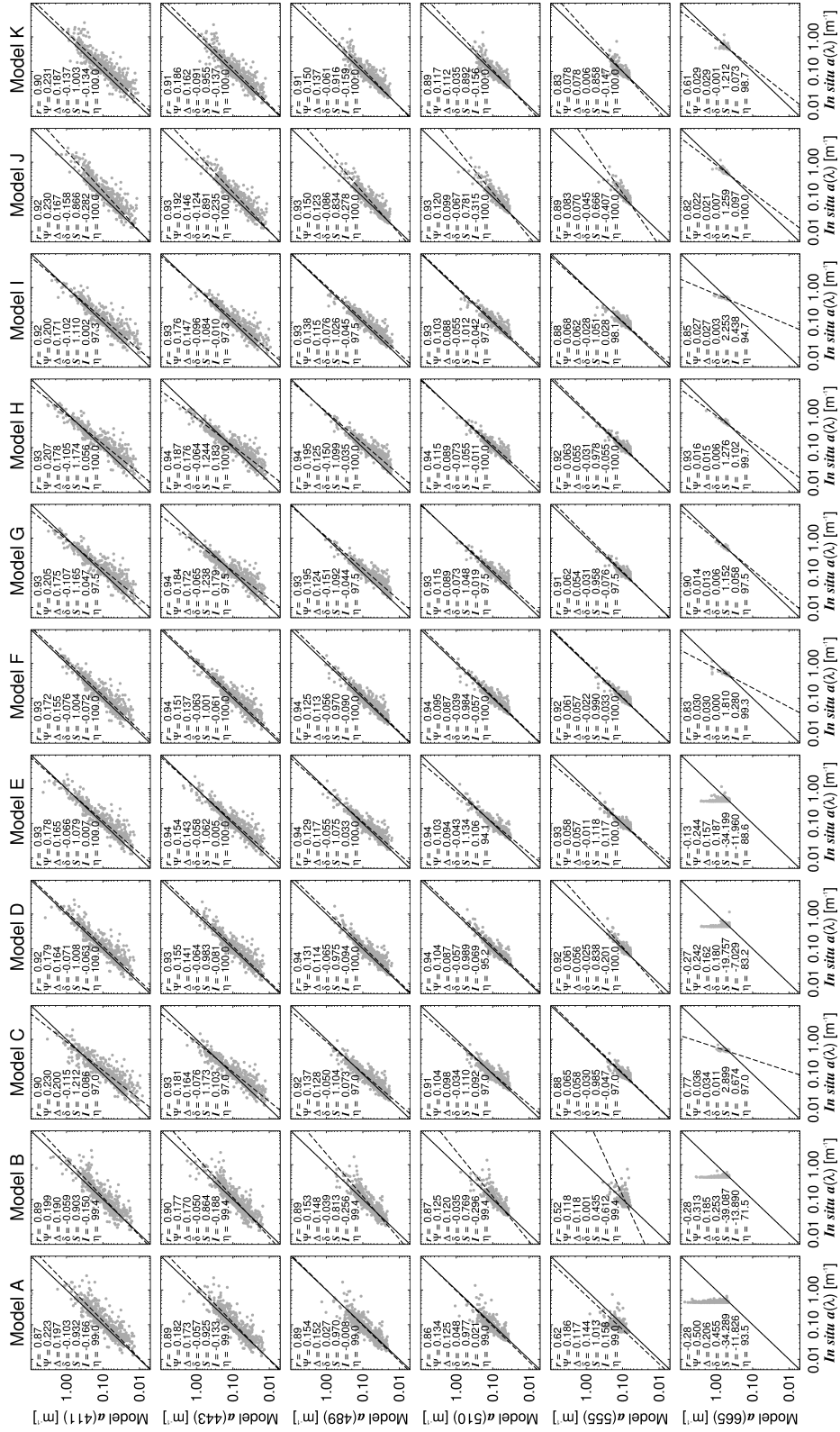


Figure 7: Scatter plots of the comparison between model and *in situ* $a(\lambda)$

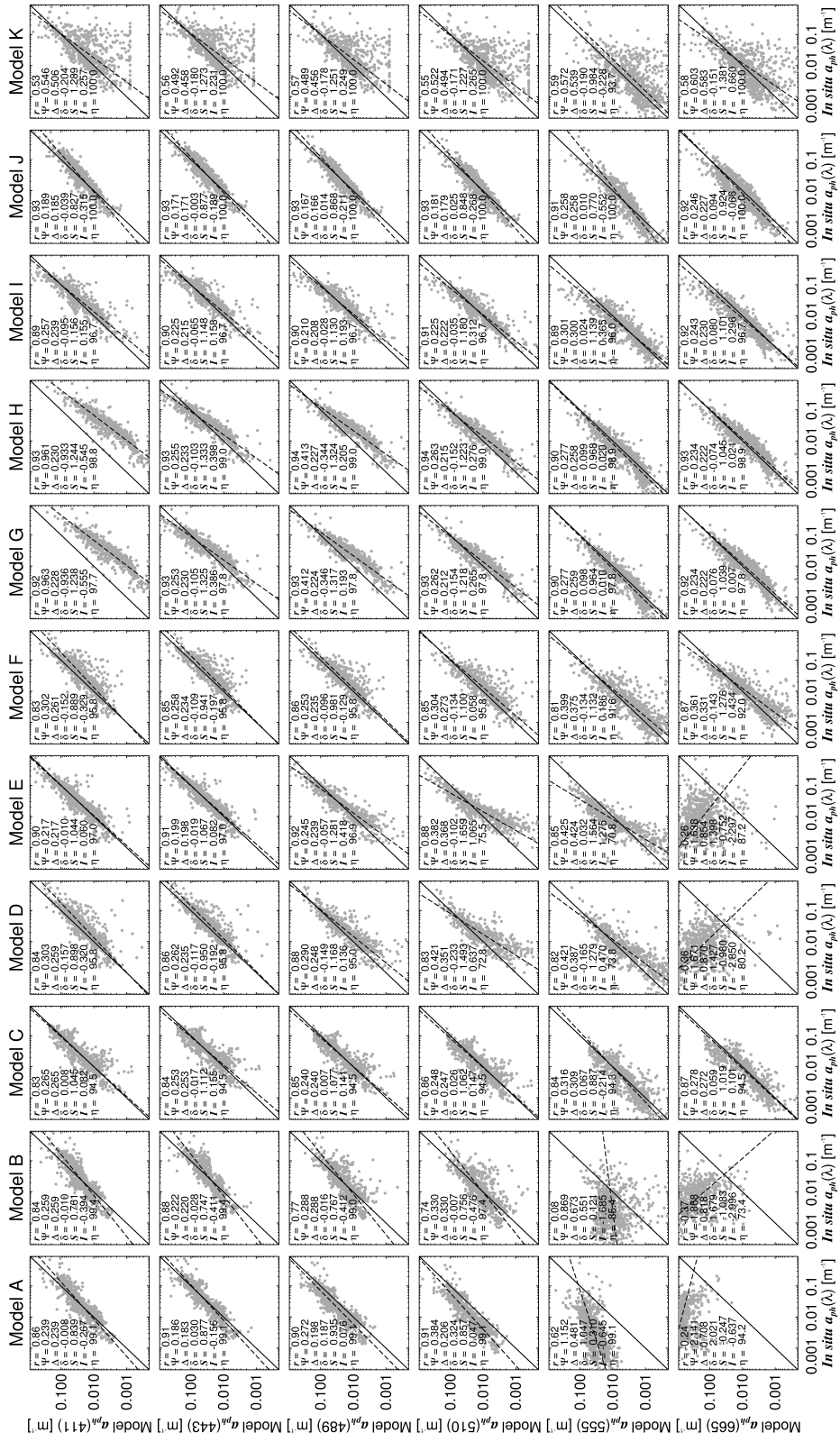


Figure 8: Scatter plots of the comparison between model and *in situ* $a_{pk}(\lambda)$

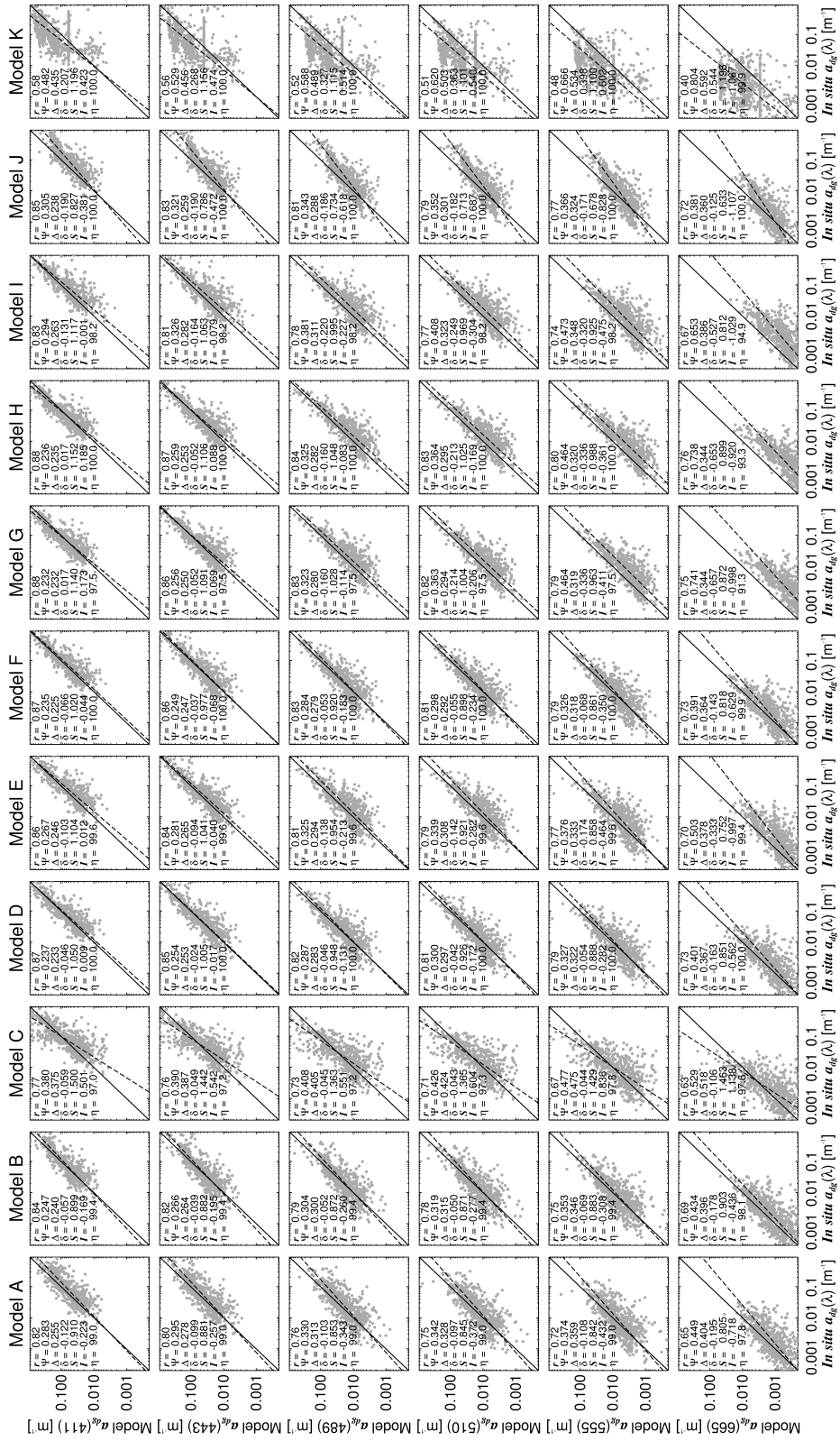


Figure 9: Scatter plots of the comparison between model and *in situ* $a_{d_g}(\lambda)$



Figure 10: Scatter plots of the comparison between model and *in situ* $b_p(\lambda)$

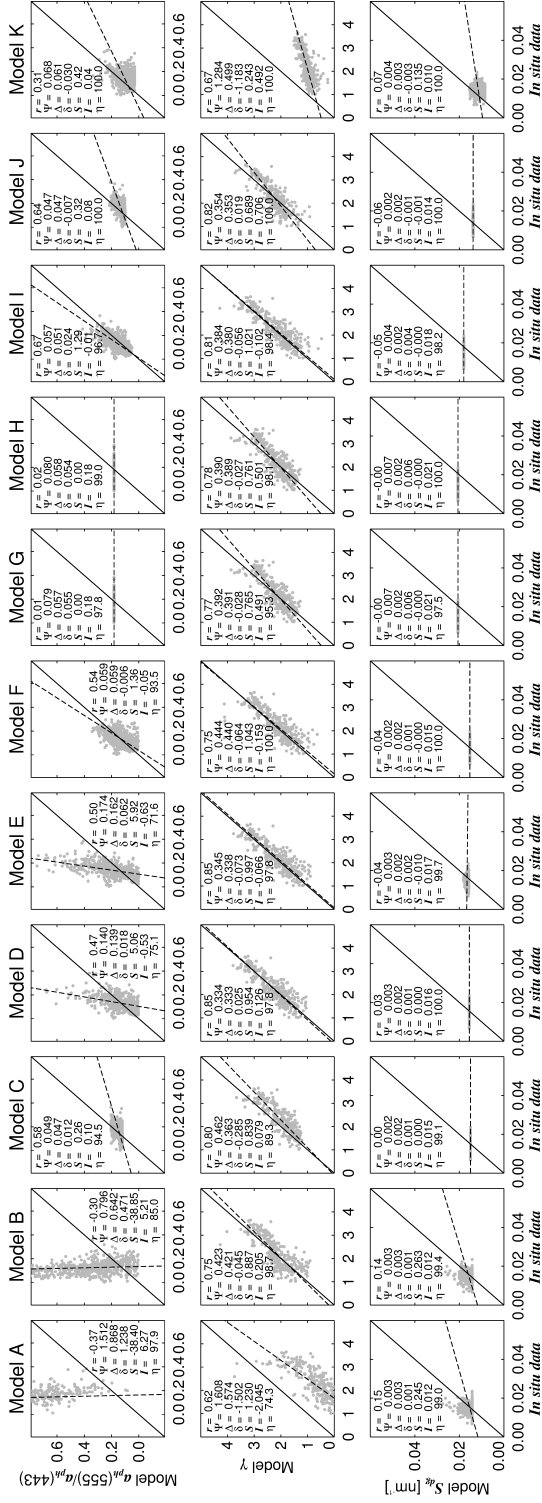


Figure 11: Scatter plots of the comparison between model and *in situ* $a_{ph}(555)/a_{ph}(443)$, γ and S_{dp} .

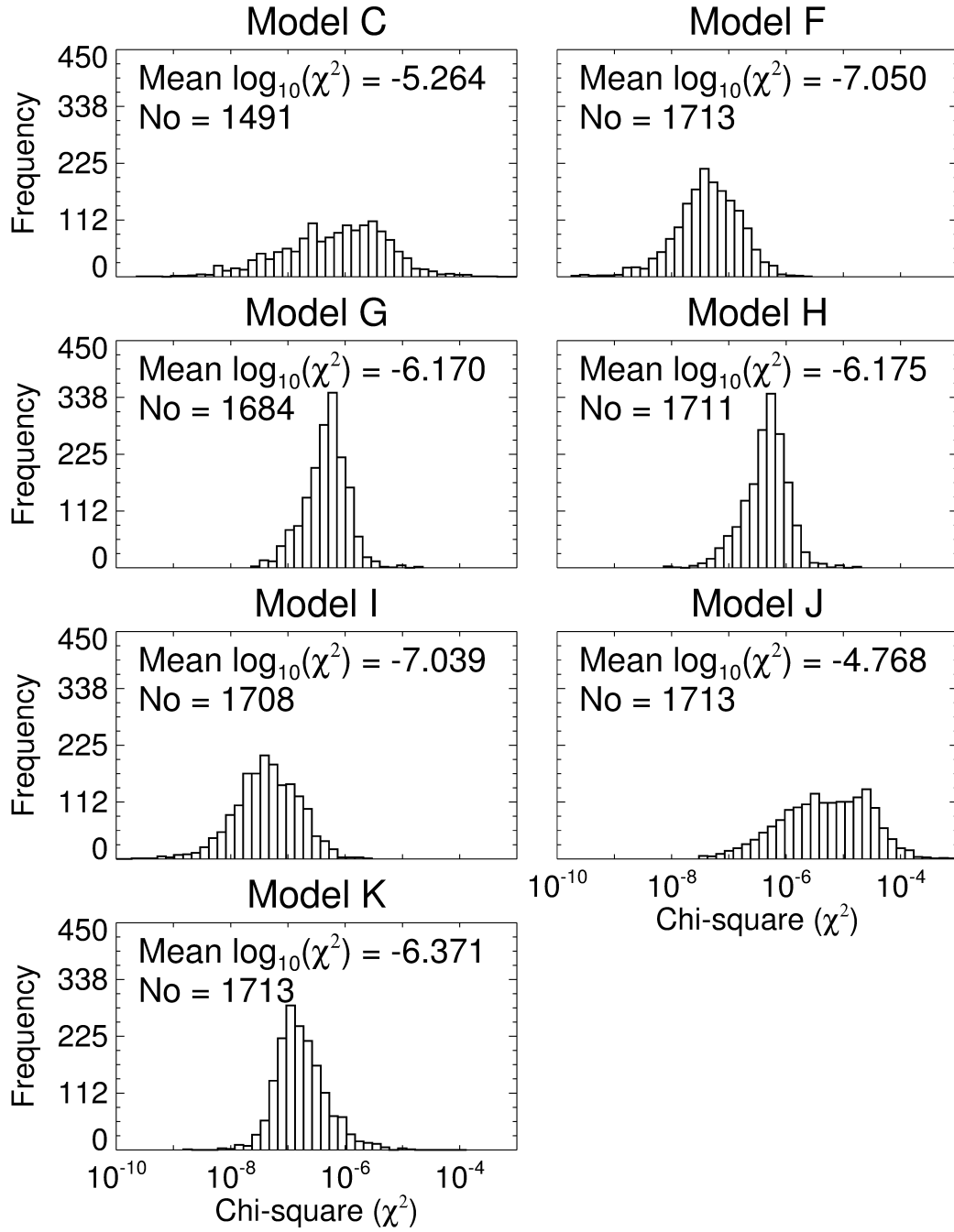


Figure 12: Results from the chi-square test.

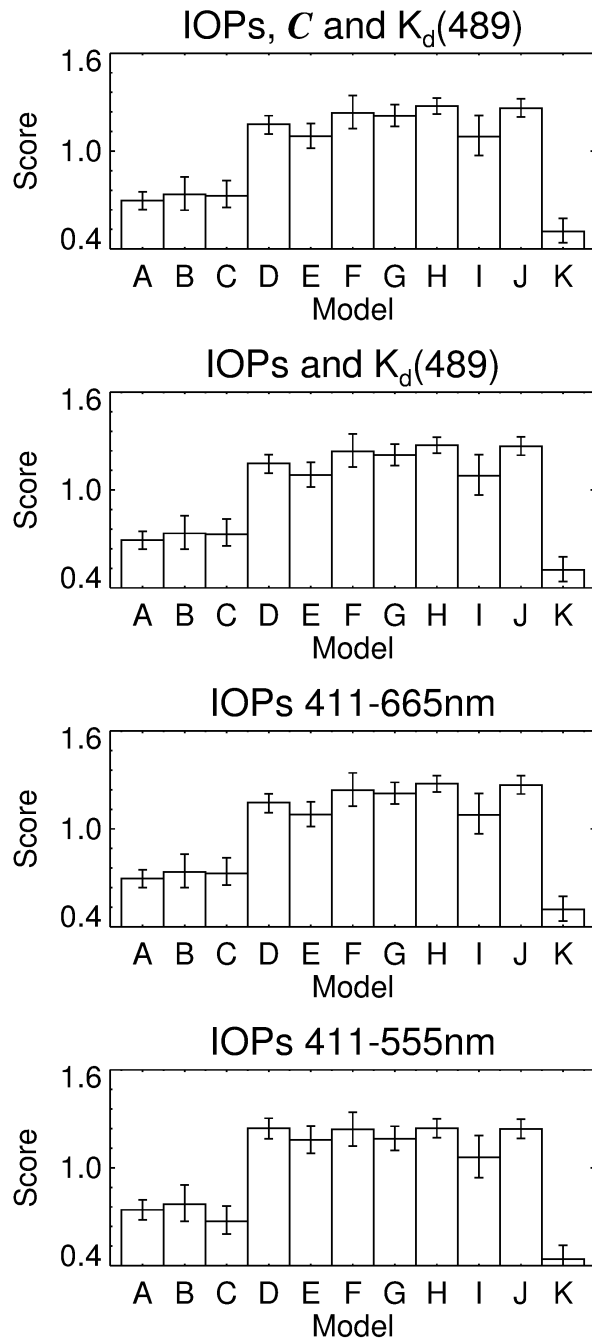


Figure 13: Results for semi-analytical models when summing all points in the classification.

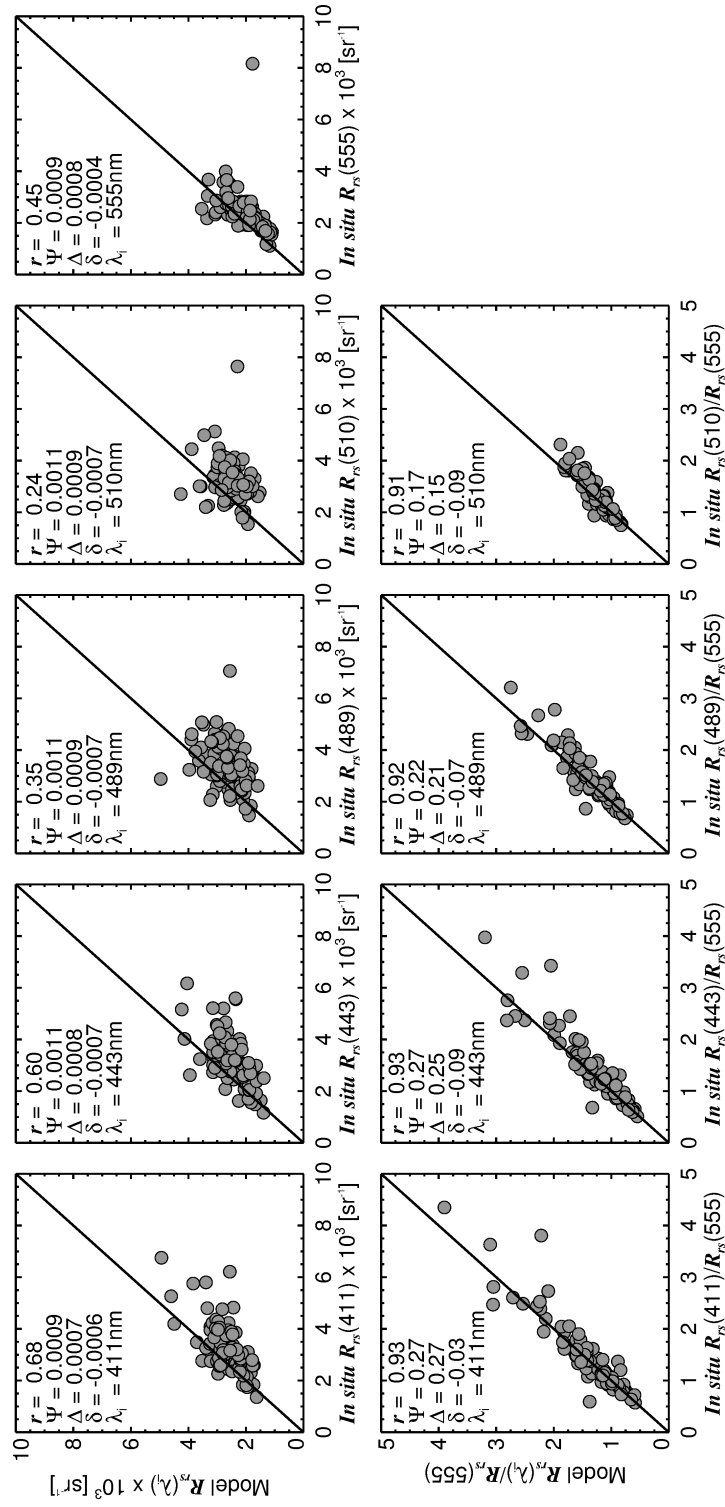


Figure 14: A comparison between measured R_{rs} and modelled R_{rs} at wavelengths from 411-555 nm for 87 samples in NOMAD with corresponding R_{rs} , a and b_b . Modelled R_{rs} has been reconstructed using *in situ* a and b_b and the approximation of Gordon et al. (1988).